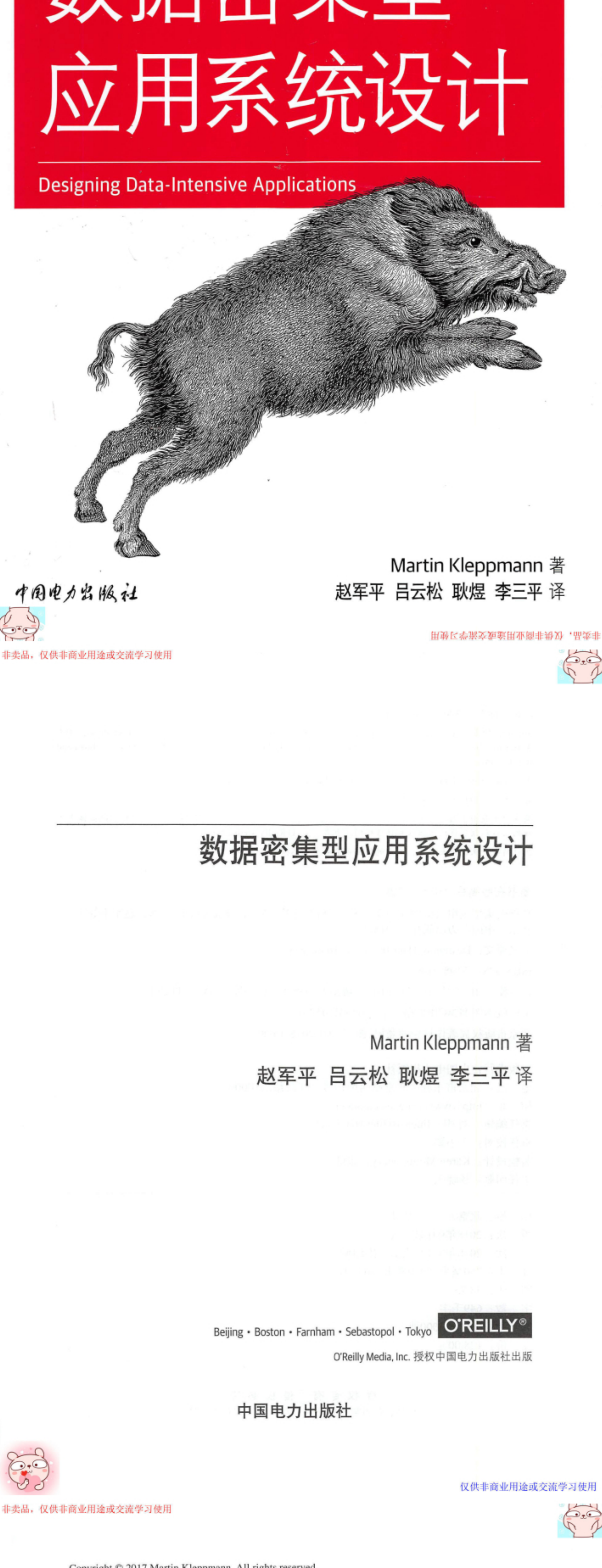




版权相关注意事项：  
1、书籍版权归著者和出版社所有  
2、本PDF来自于各个广泛的信息平台，经过整理而成  
3、本PDF仅限用于非商业用途或者个人交流研究学习使用  
4、本PDF获得者不得在互联网上以任何目的进行传播，违规者造成的法律责任和后果，违规者自负  
5、如果觉得书籍内容很赞，请一定购买正版实体书，多多支持编写高质量的图书的作者和相应的出版社！当然，如果图书内容不堪入目，质量低下，你也可以选择狠狠滴撕裂本PDF  
6、技术类书籍是拿来获取知识的，不是拿来收藏的，你得到了书籍不意味着你得到了知识，所以请不要得到书籍后就觉得沾沾自喜，要经常翻阅！！经常翻阅  
7、请于下载PDF后24小时内研究使用并删掉本PDF

版权相关注意事项：  
1、书籍版权归著者和出版社所有  
2、本PDF来自于各个广泛的信息平台，经过整理而成  
3、本PDF仅限用于非商业用途或者个人交流研究学习使用  
4、本PDF获得者不得在互联网上以任何目的进行传播，违规者造成的法律责任和后果，违规者自负  
5、如果觉得书籍内容很赞，请一定购买正版实体书，多多支持编写高质量的图书的作者和相应的出版社！当然，如果图书内容不堪入目，质量低下，你也可以选择狠狠滴撕裂本PDF  
6、技术类书籍是拿来获取知识的，不是拿来收藏的，你得到了书籍不意味着你得到了知识，所以请不要得到书籍后就觉得沾沾自喜，要经常翻阅！！经常翻阅  
7、请于下载PDF后24小时内研究使用并删掉本PDF



中国电力出版社

Martin Kleppmann 著  
赵军平 吕云松 耿煜 李三平 译

非卖品，仅供非商业用途或交流学习使用



Martin Kleppmann 著  
赵军平 吕云松 耿煜 李三平 译

Beijing • Boston • Farnham • Sebastopol • Tokyo

O'Reilly Media, Inc. 授权中国电力出版社出版

中国电力出版社

Copyright © 2017 Martin Kleppmann. All rights reserved.

Simplified Chinese Edition, jointly published by O'Reilly Media, Inc. and China Electric Power Press, 2018. Authorized translation of the English edition, 2018 O'Reilly Media, Inc., the owner of all rights to publish and sell the same.

All rights reserved including the rights of reproduction in whole or in part in any form.

英文原版由 O'Reilly Media, Inc. 出版 2017。

简体中文版由中国电力出版社出版 2018。英文原版的翻译得到 O'Reilly Media, Inc. 的授权。此简体中文版的出版和销售得到出版权和销售权的所有者——O'Reilly Media, Inc. 的许可。

版权所有，未得书面许可，本书的任何部分和全部不得以任何形式复制或。

图书在版编目（CIP）数据  
数据密集型应用系统设计 / (美) 马丁·科勒普曼 (Martin Kleppmann) 著；赵军平等译。—北京：中国电力出版社，2018.9  
书名原文：Designing Data-Intensive Applications  
ISBN 978-7-5198-2196-8  
I. ①数… II. ①马… ②赵… III. ①数据处理软件—程序设计 IV. ①TP274  
中国版本图书馆CIP数据核字(2018)第145261号  
北京市版权局著作权合同登记 图字：01-2018-1598号

出版发行：中国电力出版社  
网 址：北京市东城区北京站西街19号（邮政编码100005）  
地 址：http://www.cepp.sgcc.com.cn  
责任编辑：刘 焱 (liuchi1030@163.com)  
责任校对：王小鹏  
装帧设计：Karen Montgomery, 张健  
责任印制：杨晓东

印 刷：北京天宇星印刷厂  
版 次：2018年9月第一版  
印 次：2018年9月北京第一次印刷  
开 本：750毫米×980毫米 16开本  
印 张：33.75  
字 数：649千字  
印 数：0001—3000册  
定 价：128.00元

版 权 专 有 侵 权 必 究  
本书如有印装质量问题，我社发行部负责退换

O'Reilly Media, Inc. 介绍

O'Reilly Media通过图书、杂志、在线服务、调查研究和会议等方式传播创新知识。自1978年开始，O'Reilly一直都是前沿发展的见证者和推动者。超级极客们正在开创未来，而我们关注真正重要的技术趋势——通过放大那些“细微的信号”来刺激社会对新科技的应用。作为技术社区中活跃的参与者，O'Reilly的发展充满了对创新的倡导、创造和发扬光大。

O'Reilly为软件开发人员带来革命性的“动物书”；创建第一个商业网站（GNN）；组织了影响深远的开放源代码峰会，以至于开源软件运动以此命名；创立了Make杂志，从而成为DIY革命的主要先锋；公司一如既往地通过多种形式缔结信息与人的纽带。O'Reilly的会议和峰会集聚了众多超级极客和高瞻远瞩的商业领袖，共同描绘出开创新产业的革命性思想。作为技术人士获取信息的选择，O'Reilly现在还将先锋专家的知识传递给普通的计算机用户。无论是通过书籍出版，在线服务或者面授课程，每一项O'Reilly的产品都反映了公司不可动摇的理念——信息是激发创新的力量。

业界评论

“O'Reilly Radar博客有口皆碑。”  
——Wired

“O'Reilly凭借一系列（我希望当初我也想到了）非凡想法建立了数百万美元的业务。”  
——Business 2.0

“O'Reilly Conference是聚集关键思想领袖的绝对典范。”  
——CRN

“一本O'Reilly的书就代表一个有用、有前途、需要学习的主题。”  
——Irish Times

“Tim是位特立独行的商人，他不光放眼于最长远、最广阔的视野并且切实地按照Yogi Berra的建议去做了：‘如果你在路上遇到岔路口，走小路（岔路）’。回顾过去Tim似乎每一次都选择了小路，而且有几多次都是一闪即逝的机会，尽管大路也不错。”  
——Linux Journal

译者序

数据密集与计算密集是当今两大典型负载类型，前者以大数据为代表，后者以深度学习和HPC等为主要代表。两者各有自己的使命，也面临一些共同的挑战。其中大数据是开展深度学习的重要前提，也是当今云计算和Web服务的核心支撑技术。本书属于大数据范畴，主要采用分布式系统来处理存储数据（涉及数据库、Hadoop、NoSQL、流处理等），探讨了三个主题：系统的可靠性、可扩展性与可维护性。这三个可謂是顶级工程挑战，我们认为值得所有相关从业者重点关注。

全书深入浅出，介绍了很多简单易懂的示例。正如微软CTO以及Kafka创始人倾力推荐的那样：理论与实践相结合，层层推进，娓娓道来。知其然，更要知其所以然。每章开头都用一个地图做索引来串联技术脉络，每章结束的参考文献非常丰富，可以作为进一步学习的详细指南。可以说，本书适合分布式领域几乎所有层次的读者，包括开发者、架构师，或者相关高年级本科生、研究生等。

本书由四位译者合作完成，其中赵军平负责翻译前言、第1章和第二部分（第5章~第9章），以及全书的统筹审校；吕云松负责翻译第2章~第4章，并协助校对部分章节；耿煜负责翻译第11章和第12章，李三平负责翻译第10章。作者历时四年，倾注了很多心血，我们尽最大努力，希望不辜负使命，以使更多的读者受益。由于水平和时间有限，书中有疏漏或者不尽如人意之处，敬请广大读者批评指正。

译者

版权相关注意事项：  
1、书籍版权归著者和出版社所有  
2、本PDF来自于各个广泛的信息平台，经过整理而成  
3、本PDF仅限用于非商业用途或者个人交流研究学习使用  
4、本PDF获得者不得在互联网上以任何目的进行传播，违规者造成的法律责任和后果，违规者自负  
5、如果觉得书籍内容很赞，请一定购买正版实体书，多多支持编写高质量的图书的作者和相应的出版社！当然，如果图书内容不堪入目，质量低下，你也可以选择狠狠滴撕裂本PDF  
6、技术类书籍是拿来获取知识的，不是拿来收藏的，你得到了书籍不意味着你得到了知识，所以请不要得到书籍后就觉得沾沾自喜，要经常翻阅！！经常翻阅  
7、请于下载PDF后24小时内研究使用并删掉本PDF

非卖品，仅供非商业用途或交流学习使用

技术是推动社会进步的强大动力。数据、软件和通信等，可能会被别有用心的人所用，保护既得利益。技术需要善加利用；让弱小者的声音得到倾听，让每个人都有参与的机会，让世界免于灾难之苦。

谨以本书献给那些追逐梦想的人们。

在我看来，计算机技术是一种流行文化。流行文化只在乎认同感和参与感，活在当下，而不从关心合作、过去或者未来。我认为大部分为金钱而编写代码的人也是如此。他们不了解他们的文化来自哪里。

——Alan Kay, 接受美国 Dr. Dobbs 杂志采访（2012）

版权相关注意事项：  
1、书籍版权归著者和出版社所有  
2、本PDF来自于各个广泛的信息平台，经过整理而成  
3、本PDF仅限用于非商业用途或者个人交流研究学习使用  
4、本PDF获得者不得在互联网上以任何目的进行传播，违规者造成的法律责任和后果，违规者自负  
5、如果觉得书籍内容很赞，请一定购买正版实体书，多多支持编写高质量的图书的作者和相应的出版社！当然，如果图书内容不堪入目，质量低下，你也可以选择狠狠滴撕裂本PDF  
6、技术类书籍是拿来获取知识的，不是拿来收藏的，你得到了书籍不意味着你得到了知识，所以请不要得到书籍后就觉得沾沾自喜，要经常翻阅！！经常翻阅  
7、请于下载PDF后24小时内研究使用并删掉本PDF

非卖品，仅供非商业用途或交流学习使用





# 目录

前言 .....	1
----------	---

## 第一部分 数据系统基础

第1章 可靠、可扩展与可维护的应用系统 .....	11
---------------------------	----

认识数据系统 .....	12
可靠性 .....	14
可扩展性 .....	18
可维护性 .....	25
小结 .....	28

第2章 数据模型与查询语言 .....	33
---------------------	----

关系模型与文档模型 .....	34
数据查询语言 .....	46
图状数据模型 .....	52
小结 .....	65

第3章 数据存储与检索 .....	71
-------------------	----

数据库核心：数据结构 .....	72
事务处理与分析处理 .....	89
列式存储 .....	94
小结 .....	101







第4章 数据编码与演化 .....	109
数据编码格式 .....	110
数据流模式 .....	124
小结 .....	134

## 第二部分 分布式数据系统

第5章 数据复制 .....	145
主节点与从节点 .....	146
复制滞后问题 .....	154
多主节点复制 .....	160
无主节点复制 .....	168
小结 .....	181

第6章 数据分区 .....	189
数据分区与数据复制 .....	190
键-值数据的分区 .....	190
分区与二级索引 .....	195
分区再平衡 .....	198
请求路由 .....	202
小结 .....	204

第7章 事务 .....	211
深入理解事务 .....	212
弱隔离级别 .....	221
串行化 .....	237
小结 .....	250

第8章 分布式系统的挑战 .....	259
故障与部分失效 .....	260
不可靠的网络 .....	262
不可靠的时钟 .....	271







知识，真相与谎言 .....	282
小结 .....	292

## 第9章 一致性与共识..... 303

一致性保证 .....	304
可线性化 .....	305
顺序保证 .....	319
分布式事务与共识 .....	330
小结 .....	349

## 第三部分 派生数据

## 第10章 批处理系统..... 367

使用UNIX工具进行批处理 .....	368
MapReduce与分布式文件系统 .....	375
超越MapReduce.....	394
小结 .....	403

## 第11章 流处理系统..... 413

发送事件流.....	414
数据库与流.....	424
流处理 .....	435
小结 .....	449

## 第12章 数据系统的未来..... 461

数据集成 .....	461
分拆数据库.....	469
端到端的正确性 .....	484
做正确的事情 .....	500
小结 .....	509

## 术语表..... 521







# 前言

如果你是一位软件行业从业者，尤其是从事服务器端或者后台系统软件开发，相信近年来一定被层出不穷的商业名词所包围：NoSQL、Big Data、Web-scale、Sharding、Eventual consistency、ACID、CAP理论、云服务、MapReduce和Real-time等，所有这些其实都围绕着如何构建高效存储与数据处理这一核心主题。

过去十年，在数据库领域与分布式系统方面涌现了许多引人瞩目的进展，由此深刻地影响了如何构建上层应用系统。分析这些激动人心的变化背后，你会发现有以下几个非常重要的驱动因素：

- 互联网公司，包括Google、Yahoo!、Amazon、Facebook、LinkedIn、Microsoft，以及Twitter等，它们每天都在面对海量数据和负载，迫使其不断创新，并改进支撑系统以更有效地处理这种量级的数据。
- 商业方面因素，如敏捷开发、测试驱动和对市场机会做出快速反应等，都要求尽量缩短产品开发周期，因此系统中的数据模型也要足够灵活以方便调整。
- 免费及开源软件现在已经非常成功，在很多领域足以取代商业或者定制软件。
- 硬件方面，CPU主频增长日趋缓慢，而多核系统成为新常态，网络速度则依旧保持快速发展，这就意味着并行分布式系统将会成为业界主流。
- 如今一个不起眼的小公司，也完全有能力构建起大型分布式系统：跨机器甚至跨地域的数据中心，因为现在有了逐渐普及的IaaS云服务（例如AWS）。
- 很多服务现在都要求高可用。道理很简单，系统失效或者维护时间越长，其损失代价越大，甚至大到无法承受。







正是由于这些技术，“数据密集型应用”（Data-Intensive Applications）把很多不可能变成了可能。那么什么算是“数据密集型”（data-intensive）呢？对于一个应用系统，如果“数据”是其成败决定性因素，包括数据的规模、数据的复杂度或者数据产生与变化的速率等，我们就可以称为“数据密集型应用系统”；与之对应的是计算密集型（Compute-Intensive），CPU主频往往是后者最大的制约瓶颈。

目前已经有足够丰富的技术或者工具来辅助、帮助我们开发自己的数据密集型应用，包括存储、处理等方面，而这些技术本身也在快速演进之中。例如很多人在关注新的NoSQL系统，但实际上消息队列、缓存、搜索引擎、批处理与流处理框架等相关技术也非常重要，事实上，很多应用系统总是会集成组合上述多种技术。

开篇所罗列的那一堆商业味颇浓的名词，某种程度上也彰显了当前大时代氛围，机会多自是好事。然而，作为一名软件工程师或者架构师，仍需秉持严谨的态度，深入理解繁杂词汇背后系统设计所面临的优劣权衡，只有这样才能为我所用，构建好自己的系统。所以，抛开那些商业名词，我们要的是深入的探索。

幸好，软件千变万化，终有若干理念贯穿其中。无论你用的什么数据系统，如果可以掌握背后的设计理念，何种工具适用于何种场景，如何最佳使用，又有哪些陷阱（坑），诸如此类，自然会胸有成竹，而这也是本书写作的初衷。

所以，本书旨在帮助大家更好地驾驭处理数据和存储数据相关技术。它不是针对某个特定软件的介绍手册，也不是纯理论的习题。我们会深入探讨一些成功的数据系统案例，剖析其中的技术要点；或许在很多流行的分布式系统里都有它们的身影，正是这些关键技术有效应对了许多生产环境对扩展性、性能和可靠性的苛刻要求。

我们将对这些典型系统深入展开讨论，梳理其核心算法，探讨其设计理念和背后的权衡之道。在此过程中，尝试总结某些经验法则来重新审视系统架构。了解怎么工作自然重要，但更重要的是要思考它为什么这样工作，正所谓“知其然，知其所以然”。

读完本书之后，你应该会对哪些技术适用于哪些场景，常见工具如何搭配来构建应用系统等有所得。或许你不会马上就决定动手开发一个全新的数据库引擎（这几乎完全不需要），但是，关于系统的本质，相信你一定有新的认识与判断力：系统行为是否合理，架构设计如果权衡，个中症状如何处理等，将会更加游刃有余。

## 本书适合哪些读者？

首先你的应用如果包含服务器端、后台逻辑来存储、处理数据，或者你的应用需要联网，例如Web程序、移动程序以及联网的传感器等，这本书将非常适合。







本书主要针对软件工程师、软件架构师以及技术经理等，特别是那些需要对系统架构做出权衡决定的人，例如需要选择一些工具和软件来解决特定问题，或者如何最佳适用这些现有工具。退一步，如果不需要做这些决定，本书也可以帮助你更好地理解这些技术的优缺点。

当然你最好有一些Web程序或者网络程序的经验，了解一些基本的关系型数据和SQL；虽然不是必须，但如果通晓NoSQL或相关系统，那再好不过。如果有常见网络协议如TCP、HTTP等基本知识也会帮助很大。选择何种编程语言或者框架对于阅读本书没有太大影响。

如果以下若干条适用于你，可能会从本书中有所获益：

- 需要学习了解系统扩展性方面技术，例如支撑百万用户级的Web和移动程序。
- 需要构建高可用（减少宕机影响）和健壮运行的系统。
- 需要有效方法来提高长时间运行系统的可维护性，应对规模增长、技术和需求不断发生变化等。
- 对系统如何工作有种天然的兴趣或者探索精神，特别是大型Web和在线系统。本书会例举若干典型数据库和数据处理、分析系统，对其设计亮点逐一展开分析，也是乐事一件。

当我在介绍可扩展的系统时，有人会问：“我又不是Google或者Amazon，这些高大上的扩展性太遥远，老老实实用关系数据库就得了。”在特定语境下这的确有道理，打造一个根本不需要的扩展性系统是在浪费精力，而且会丧失其他方面的灵活性，这其实是种过早的优化。然而我要说，选择何种工具也很有讲究或者说极其重要，不同的技术各有优劣，事实上正如本书将要揭示的，关系数据库很重要，但它并非解决数据问题的终极方案。

## 本书涵盖范围

本书不会就如何安装、配置特定软件包或者API等做太多细节介绍，这方面已有足够多的文档。我们会着重探讨不同的设计理念以及如何权衡，这些通常才是系统的核心，也会例举实践中不同的系统产品最终采用了哪些不同的理念。

本书的电子版已包含所有在线资源的链接，这些链接在书稿付梓之时都做了验证，然而大家都能理解，这些网络链接可能在一段时间之后会失效。如果不巧碰到了失效的链接，或者阅读的是纸质版本，可以随时用搜索引擎来查询相关资料。对于学术论







文，推荐用Google Scholar搜索文献标题来下载PDF版本，或者随时访问<https://github.com/lept/ddia-references>，我们维护了所有最新的链接。

本书主要关注数据处理系统架构方面以及如何集成到数据密集型应用系统中。而关于部署、运营、安全、管理等相关方面，当然也非常重要和足够复杂，不过它们并非本书的重点，也无意做些肤浅的、聊胜于无的介绍，我们认为它们值得单独成书。

书中所涵盖的诸多技术总体都可以纳入到“大数据”的范畴。然而，大数据这个商业词汇有些过于滥用或者太模糊，不太适合在严谨的工程领域展开讨论。本书倾向于使用更加明确的用语，例如单节点、分布式系统，或者在线/交互式分析以及离线/批处理系统等。

本书更偏爱免费及开源软件（Free and Open Source Software, FOSS），主要是通过源码阅读、必要的修改，然后实地执行可以帮助理解系统究竟是如何工作的，此外开放系统也能减少对特定供应商的锁定问题。当然，我们并非局限于此，也会讨论相关商业软件，这包括闭源、软件即服务（Software as a Service）或者一些只出现在文献里但并未公开发行的一些公司内部软件等。

## 本书内容安排

全书分为三大部分：

1. 第一部分，主要讨论有关增强数据密集型应用系统所需的若干基本原则。首先开篇第1章即瞄准目标：可靠性、可扩展性与可维护性，如何认识这些问题以及如何达成目标。第2章我们比较了多种不同的数据模型和查询语言，讨论各自的适用场景。接下来第3章主要针对存储引擎，即数据库是如何安排磁盘结构从而提高检索效率。第4章转向数据编码（序列化）方面，包括常见模式的演化历程。
2. 第二部分，我们将从单机的数据存储转向跨机器的分布式系统，这是扩展性的重要一步，但随之而来的是各种挑战。所以将依次讨论数据远程复制（第5章）、数据分区（第6章）以及事务（第7章）。接下来的第8章包括分布式系统的更多细节，以及分布式环境如何达成一致性与共识（第9章）。
3. 第三部分，主要针对产生派生数据的系统，所谓派生数据主要指在异构系统中，如果无法用一个数据源来解决所有问题，那么一种自然的方式就是集成多个不同的数据库、缓存模块以及索引模块等。首先第10章以批处理开始来处理派生数据，紧接着第11章采用流式处理。第12章总结之前介绍的多种技术，并分析讨论未来构建可靠、可扩展和可维护应用系统可能的方向或方法。







## 参考资料与进一步阅读

书中所讨论的很多内容可能在其他地方以某种方式有所呈现，包括会议演讲、学术论文、博客、源代码、Bug追踪系统、邮件列表，或是工程师间的交流。本书试图汇集其中最重要的精华呈现给读者，同时也附上相关原始出处的链接。每章最后列出的参考索引是一个很好的资源，可以帮助更深入地去了解某一方面，而且这些内容绝大多数都可以在线访问。

## O'Reilly Safari

Safari（之前的Safari Books Online）是一个针对企业、政府、教育工作者以及个人的会员制培训和参考平台。

订阅者可以访问来自超过250个出版社的大量书籍、培训视频、学习路径、交互式教程，以及精选的播放列表，这些出版社包括 O'Reilly Media、Harvard Business Review、Prentice Hall Professional、Addison-Wesley Professional、Microsoft Press、Sams、Que、Peachpit Press、Adobe、Focal Press、Cisco Press、John Wiley & Sons、Syngress、Morgan Kaufmann、IBM Redbooks、Packt、AdobePress、FT Press、Apress、Manning、New Riders、McGraw-Hill、Jones & Bartlett等。

更多信息请访问<http://orieilly.com/safari>。

## 联系我们

请将关于本书的意见和问题通过以下地址提供给出版商：

美国：

O'Reilly Media, Inc.  
1005 Gravenstein Highway North  
Sebastopol, CA 95472

中国：

北京市西城区西直门南大街2号成铭大厦C座807室（100035）  
奥莱利技术咨询（北京）有限公司

针对这本书，我们还建有一个网页，列出了有关勘误、示例和其他信息。可以通过以下地址访问这个页面：<http://bit.ly/designing-data-intensive-apps>。







有关本书技术方面的问题和评论，请联系[bookquestions@oreilly.com](mailto:bookquestions@oreilly.com)。

若想了解O'Reilly图书、培训课程、会议和新闻的更多信息，请访问我们的网站，地址是：<http://www.oreilly.com>，<http://www.oreilly.com.cn>。

我们的Facebook: <http://facebook.com/oreilly>。

我们的Twitter: <http://twitter.com/oreillymedia>。

我们的YouTube: <http://www.youtube.com/oreillymedia>。

## 致谢

本书融合了学术研究和工程实践的很多经验，是大量他人思想、知识的融合与系统化。在计算机领域，人们往往会被新的东西所吸引，但我认为前人有很多优秀的工作和积累依然值得深入学习。本书有超过800篇参考文献，包括论文、博客文章、演讲、产品文档等，对我来说他们都是本书非常宝贵的学习资源。在此，我衷心感谢原作者们无私的知识分享。

我也从很多交流对话中领悟颇多，感谢这些花费了大量时间与我讨论想法、耐心解释的朋友们。特别是要感谢Joe Adler、Ross Anderson、Peter Bailis、Márton Balassi、Alastair Beresford、Mark Callaghan、Mat Clayton、Patrick Collison、Sean Cribbs、Shirshanka Das、Niklas Ekström、Stephan Ewen、Alan Fekete、Gyula Fóra、Camille Fournier、Andres Freund、John Garbutt、Seth Gilbert、Tom Haggett、Pat Helland、Joe Hellerstein、Jakob Homan、Heidi Howard、John Hugg、Julian Hyde、Conrad Irwin、Evan Jones、Flavio Junqueira、Jessica Kerr、Kyle Kingsbury、Jay Kreps、Carl Lerche、Nicolas Liochon、Steve Loughran、Lee Mallabone、Nathan Marz、Caitie McCaffrey、Josie McLellan、Christopher Meiklejohn、Ian Meyers、Neha Narkhede、Neha Narula、Cathy O'Neil、Onora O'Neil、Ludovic Orban、Zoran Perkov、Julia Powles、Chris Riccomini、Henry Robinson、David Rosenthal、Jennifer Rullmann、Matthew Sackman、Martin Scholl、Amit Sela、Gwen Shapira、Greg Spurrier、Sam Stokes、Ben Stopford、Tom Stuart、Diana Vasile、Rahul Vohra、Pete Warden和Brett Wooldridge。

还有些朋友认真阅读了本书草稿并提供了反馈，这对本书的最终成形帮助很大，在此也衷心感谢，他们是Raul Agepati、Tyler Akidau、Mattias Andersson、Sasha







Baranov、Veena Basavaraj、David Beyer、Jim Brikman、Paul Carey、Raul Castro Fernandez、Joseph Chow、Derek Elkins、Sam Elliott、Alexander Gallego、Mark Grover、Stu Halloway、Heidi Howard、Nicola Kleppmann、Stefan Kruppa、Bjorn Madsen、Sander Mak、Stefan Podkowinski、Phil Potter、Hamid Ramazani、Sam Stokes和Ben Summers。当然，对于本书的任何错误或者问题，由我本人承担全部责任。

非常感谢编辑Marie Beaugureau、Mike Loukides、Ann Spencer和O'Reilly的所有团队，他们忍耐着我缓慢的写作进度以及各种要求，为这本书付出了很多心血。同时感谢Rachel Head，为本书润色不少。我要感谢Alastair Beresford、Susan Goodhue、Neha Narkhede和Kevin Scott，尽管我还有其他的工作安排，但他们给了我非常大的写作时间和自由。

特别感谢Shabbir Diwan和Edie Freedman，他们为每章配上了细致的前导地图。这个想法非常棒，用创造性的地图方式来展示内容脉络，最终效果非常漂亮和酷炫！

最后，感谢我的家人与朋友，没有他们的支持，我肯定无法熬过这漫长的四年写作过程。你们是我的挚爱！







## 第一部分

# 数据系统基础

本书前4章总结了适用于所有数据系统的基本思想，既包括单机运行环境，也包括分布式集群环境：

1. 第1章介绍相关术语与方法，这些术语等将贯穿于全书。例如重点关注的可靠性、可扩展性与可维护性设计目标，以及达到这些目标的基本方法。
2. 第2章对比多种不同的数据模型与查询语言，从开发者角度来看，这些是不同数据库系统最显著的区别，我们也会讨论不同模型的具体适用场景。
3. 第3章深入数据库系统内部的核心存储引擎，详细解析数据库如何设计磁盘布局。针对不同的工作负载如何优化其引擎，而正确的设计选型对系统将产生巨大的影响。
4. 第4章比较不同的数据编码格式和序列化技术，特别是当上层应用需求多变而模型也需要灵活调整时，该如何最佳使用这些技术。

之后，在第二部分，我们将转向分布式场景下的一些典型问题和解决方案。

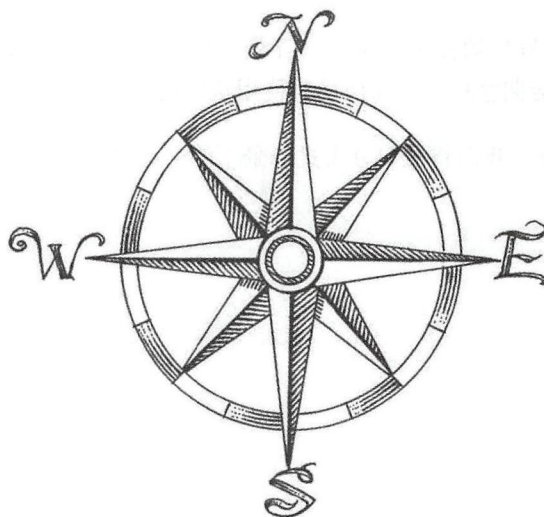






代码一景

数据密集型应用系统







## 第1章

# 可靠、可扩展与可维护的应用系统

互联网做得太出色了，以至于很多人把它看作某种像太平洋一样的自然资源，而非人造的。你还记得上一个达到如此规模而又这样健壮的技术是什么？

——Alan Kay, Dr. Dobbs杂志采访 (2012)

当今许多新型应用都属于数据密集型（data-intensive），而不是计算密集型（compute-intensive）。对于这些类型应用，CPU的处理能力往往不是第一限制性因素，关键在于数据量、数据的复杂度及数据的快速多变性。

数据密集型应用通常也是基于标准模块构建而成，每个模块负责单一的常用功能。例如，许多应用系统都包含以下模块：

- 数据库：用以存储数据，这样之后应用可以再次访问。
- 高速缓存：缓存那些复杂或操作代价昂贵的结果，以加快下一次访问。
- 索引：用户可以按关键字搜索数据并支持各种过滤。
- 流式处理：持续发送消息至另一个进程，处理采用异步方式。
- 批处理：定期处理大量的累积数据。

这些模块也许看上去习以为常，主要是因为这些“数据处理系统”（data systems）







已经做了非常漂亮的抽象，以至于我们已经习惯于拿来即用，而没有做太多深入地思考。由于目前业界已有很多数据方案可供选择，当需要构建一个新应用时，相信大多数开发者不会从头开始，例如写一个全新的数据存储引擎。

但情况并非如此简单。确实已有很多数据库系统，但因为需求和设计目标的差异，个中精妙都不尽相同。缓存和索引方案与之类似。因此在构建某个特定应用时，我们总是需要弄清楚哪些组件、哪些方法最适合自己，并且当单个组件无法满足需求而必须组合使用时，总要面临更多的技术挑战。

本书既介绍系统原理，也注重具体实践，通过二者结合的章节之旅慢慢展示如何构建数据密集型应用。我们还将探索现有工具的相同之处、差异所在，以及它们如何达成预设目标。

在本章，我们将首先探讨所关注的核心设计目标：可靠、可扩展与可维护的数据系统。澄清本源，解析处理之道，建立后续章节所需的基本要点。在接下来的章节中，我们将层层推进，深入探讨系统设计时所面临的种种抉择与权衡。

## 认识数据系统

我们通常将数据库、队列、高速缓存等视为不同类型的系统。虽然数据库和消息队列存在某些相似性，例如两者都会保存数据（至少一段时间），但他们却有着截然不同的访问模式，这就意味着不同的性能特征和设计实现。

那么为什么本书将它们归为一大类即“数据系统”（data system）呢？

首先，近年来出现了许多用于数据存储和处理的新工具。它们针对各种不同的应用场景进行优化，不适合再归为传统类型<sup>[1]</sup>。例如，Redis既可以用于数据存储也适用于消息队列，Apache Kafka作为消息队列也具备了持久化存储保证。系统之间的界限正在变得模糊。

其次，越来越多的应用系统需求广泛，单个组件往往无法满足所有数据处理与存储需求。因而需要将任务分解，每个组件负责高效完成其中一部分，多个组件依靠应用层代码驱动有机衔接起来。

举个例子，假定某个应用包含缓存层（例如Memcached）与全文索引服务器（如Elasticsearch或Solr），二者与主数据库保持关联，通常由应用代码负责缓存、索引与主数据库之间的同步，如图1-1所示（具体技术将在后面的章节中详细介绍）。





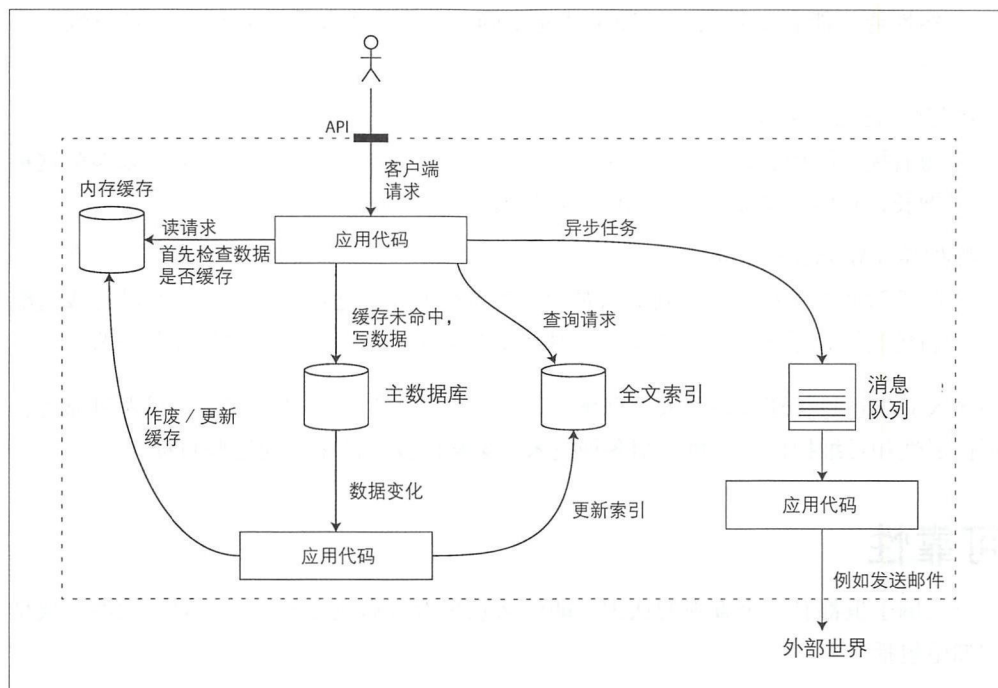


图1-1：一种数据系统架构，它包含了多个不同组件

在上面例子中，组合使用了多个组件来提供服务，而对外提供服务的界面或者API会隐藏很多内部实现细节。这样基本上我们基于一个个较小的、通用的组件，构建而成一个全新的、专用的数据系统。这样的集成数据系统会提供某些技术保证，例如，缓存要正确刷新以保证外部客户端看到一致的结果。现在可以说，你既是一名应用开发者，也是一名数据系统设计师。

设计数据系统或数据服务时，一定会碰到很多棘手的问题。例如，当系统内出现了局部失效时，如何确保数据的正确性与完整性？当发生系统降级（degrade）时，该如何为客户提供一致的良好表现？负载增加时，系统如何扩展？友好的服务API该如何设计？

影响数据系统设计的因素有很多，其中包括相关人员技能和经验水平、遗留系统依赖性、交付周期、对不同风险因素的容忍度、监管合规等。这些因素往往因时因地而异。本书将专注于对大多数软件系统都极为重要的三个问题：

### 可靠性 (Reliability)

当出现意外情况如硬件、软件故障、人为失误等，系统应可以继续正常运转：虽







然性能可能有所降低，但确保功能正确。具体请参阅本章后面的“可靠性”一节。

### 可扩展性 (Scalability)

随着规模的增长，例如数据量、流量或复杂性，系统应以合理的方式来匹配这种增长，具体请参阅本章后面的“可扩展性”一节。

### 可维护性 (Maintainability)

随着时间的推移，许多新的人员参与到系统开发和运维，以维护现有功能或适配新场景等，系统都应高效运转。具体请参阅本章后面的“可维护性”一节。

很多人对上述目标还欠缺深入的理解。追本溯源，本章后续部分将深入思考可靠性、可扩展性和可维护性，并将介绍各种技术、架构以及算法来实现这些目标。

## 可靠性

每个人脑子里都有一个直观的认识，即什么意味着可靠或者不可靠。对于软件，典型的期望包括：

- 应用程序执行用户所期望的功能。
- 可以容忍用户出现错误或者不正确的软件使用方法。
- 性能可以应对典型场景、合理负载压力和数据量。
- 系统可防止任何未经授权的访问和滥用。

如果所有上述目标都要支持才算“正常工作”，那么我们可以认为可靠性大致意味着：即使发生了某些错误，系统仍可以继续正常工作。

可能出错的事情称为错误 (faults) 或故障，系统可应对错误则称为容错 (fault-tolerant) 或者弹性 (resilient)。前一个词略显误导：似乎暗示着系统可以容忍各种可能的故障类型，显然实际中这是不可能的。举一个夸张一些的例子，如果整个地球（及其上的所有服务器）都被黑洞吞噬，那么要在这个级别容错就意味着必须在宇宙范围内进行系统冗余。试想，这将是天价的预算。因此，容错总是指特定类型的故障，这样的系统才更有实际意义。

注意，故障与失效 (failure) 不完全一致<sup>[2]</sup>。故障通常被定义为组件偏离其正常规格，而失效意味系统作为一个整体停止，无法向用户提供所需的服务。我们不太可能将故障概率降低到零，因此通常设计容错机制来避免从故障引发系统失效。本书将介绍在不可靠组件基础上构建可靠性系统的相关技术。





在这种容错系统中，用于测试目的，可以故意提高故障发生概率，例如通过随机杀死某个进程，来确保系统仍保持健壮。很多关键的bug实际上正是由于错误处理不当而造成的<sup>[3]</sup>。通过这种故意引发故障的方式，来持续检验、测试系统的容错机制，增加对真实发生故障时应对的信心。Netflix的Chaos Monkey系统<sup>[4]</sup>就是这种测试的典型例子。

虽然我们通常倾向于容忍故障而不是预防故障，但是也存在“预防胜于治疗”的情况，安全问题就是一例，例如：如果攻击者破坏了系统并窃取了敏感数据，则该事件造成的影响显然无法被撤销。然而，本书主要针对那些影响可以被消除的故障类型，接下来详细介绍。

## 硬件故障

当我们考虑系统故障时，对于硬件故障总是很容易想到：硬盘崩溃，内存故障，电网停电，甚至有人误拔掉了网线。任何与大型数据中心合作过的人都可以告诉你，当有很多机器时，这类事情迟早会发生。

有研究证明硬盘的平均无故障时间（MTTF）约为10~50年<sup>[5,6]</sup>。因此，在一个包括10 000个磁盘的存储集群中，我们应该预期平均每天有一个磁盘发生故障。

我们的第一个反应通常是为硬件添加冗余来减少系统故障率。例如对磁盘配置RAID，服务器配备双电源，甚至热插拔CPU，数据中心添加备用电源、发电机等。当一个组件发生故障，冗余组件可以快速接管，之后再更换失效的组件。这种方法可能并不能完全防止硬件故障所引发的失效，但还是被普遍采用，且在实际中也确实可以让系统不间断运行长达数年。

直到最近，采用硬件冗余方案对于大多数应用场景还是足够的，它使得单台机器完全失效的概率降为非常低的水平。只要可以将备份迅速恢复到新机器上，故障的停机时间在大多数应用中并不是灾难性的。而多机冗余则只对少量的关键应用更有意义，对于这些应用，高可用性是绝对必要的。

但是，随着数据量和应用计算需求的增加，更多的应用可以运行在大规模机器之上，随之而来的硬件故障率呈线性增长。例如，对于某些云平台（如Amazon Web Services, AWS），由于系统强调的是总体灵活性与弹性<sup>注1</sup>而非单台机器的可靠性，虚拟机实例经常会事先无告警的情况下出现无法访问问题<sup>[7]</sup>。

---

注1：具体参阅本章后面的“应对负载增加的方法”。







因此，通过软件容错的方式来容忍多机失效成为新的手段，或者至少成为硬件容错的有力补充。这样的系统更具有操作便利性，例如当需要重启计算机时为操作系统打安全补丁，可以每次给一个节点打补丁然后重启，而不需要同时下线整个系统（即滚动升级，详见第4章）。

## 软件错误

我们通常认为硬件故障之间多是相互独立的：一台机器的磁盘出现故障并不意味着另一台机器的磁盘也要失效。除非存在某种弱相关（例如一些共性原因，如服务器机架中的温度过高），否则通常不太可能出现大量硬件组件同时失效的情况。

另一类故障则是系统内的软件问题<sup>[8]</sup>。这些故障事先更加难以预料，而且因为节点之间是由软件关联的，因而往往会导致更多的系统故障<sup>[5]</sup>。例如：

- 由于软件错误，导致当输入特定值时应用服务器总是崩溃。例如，2012年6月30日发生闰秒，由于Linux内核中的一个bug，导致了很多应用程序在该时刻发生挂起<sup>[9]</sup>。
- 一个应用进程使用了某些共享资源如CPU、内存、磁盘或网络带宽，但却不幸失控跑飞了。
- 系统依赖于某些服务，但该服务突然变慢，甚至无响应或者开始返回异常的响应。
- 级联故障，其中某个组件的小故障触发另一个组件故障，进而引发更多的系统问题<sup>[10]</sup>。

导致软件故障的bug通常会长时间处于引而不发的状态，直到碰到特定的触发条件。这也意味着系统软件其实对使用环境存在某种假设，而这种假设多数情况都可以满足，但是在特定情况下，假设条件变得不再成立<sup>[11]</sup>。

软件系统问题有时没有快速解决办法，而只能仔细考虑很多细节，包括认真检查依赖的假设条件与系统之间交互，进行全面的测试，进程隔离，允许进程崩溃并自动重启，反复评估，监控并分析生产环节的行为表现等。如果系统提供某些保证，例如，在消息队列中，输出消息的数量应等于输入消息的数量，则可以不断地检查确认，如发现差异则立即告警<sup>[12]</sup>。

## 人为失误

设计和构建软件系统总是由人类完成，也是由人来运维这些系统。即使有时意图是





好的，但人却无法做到万无一失。例如，一项针对大型互联网服务的调查发现，运维者的配置错误居然是系统下线的首要原因，而硬件问题（服务器或网络）仅在10%~25%的故障中有所影响<sup>[13]</sup>。

如果我们假定人是不可靠的，那么该如何保证系统的可靠性呢？可以尝试结合以下多种方法：

- 以最小出错的方式来设计系统。例如，精心设计的抽象层、API以及管理界面，使“做正确的事情”很轻松，但搞坏很复杂。但是，如果限制过多，人们就会想法来绕过它，这会抵消其正面作用。因此解决之道在于很好的平衡。
- 想办法分离最容易出错的地方、容易引发故障的接口。特别是，提供一个功能齐全但非生产用的沙箱环境，使人们可以放心的尝试、体验，包括导入真实的数据，万一出现问题，不会影响真实用户。
- 充分的测试：从各单元测试到全系统集成测试以及手动测试<sup>[3]</sup>。自动化测试已被广泛使用，对于覆盖正常操作中很少出现的边界条件等尤为重要。
- 当出现人为失误时，提供快速的恢复机制以尽量减少故障影响。例如，快速回滚配置改动，滚动发布新代码（这样任何意外的错误仅会影响一小部分用户），并提供校验数据的工具（防止旧的计算方式不正确）。
- 设置详细而清晰的监控子系统，包括性能指标和错误率。在其他行业称为遥测（Telemetry），一旦火箭离开地面，遥测对于跟踪运行和了解故障至关重要<sup>[14]</sup>。监控可以向我们发送告警信号，并检查是否存在假设不成立或违反约束条件等。这些检测指标对于诊断问题也非常有用。
- 推行管理流程并加以培训。这非常重要而且比较复杂，具体内容已超出本书范围。

## 可靠性的重要性

可靠性绝不仅仅针对的是核电站和空中交管软件之类的系统，很多应用都需要可靠工作。商业软件中的错误会导致效率下降（如数据报告错误，甚至带来法律风险），电子商务网站的暂停会对营收和声誉带来巨大损失。

即使在所谓“非关键”应用中，我们也应秉持对用户负责的态度。例如一对父母，将其所有的照片以及他们孩子的视频存放在你的照片应用中<sup>[15]</sup>。如果不幸发生了数据库损坏，他们的感受可想而知，他们是否知道该如何从备份数据来执行恢复？







当然，也会存在其他一些情况，例如面对不太确定的市场开发原型系统，或者服务的利润微薄，有时也会牺牲一些可靠性来降低开发成本或者运营开销，对此，我们总是建议务必三思后行。

## 可扩展性

即使系统现在工作可靠，并不意味着它将来一定能够可靠运转。发生退化的一个常见原因是负载增加：例如也许并发用户从最初的10 000个增长到100 000个，或从100万到1000万；又或者系统目前要处理的数据量超出之前很多倍。

可扩展性是用来描述系统应对负载增加能力的术语。但是请注意，它并不是衡量一个系统的一维指标，谈论“X是可扩展”或“Y不扩展”没有太大意义。相反，讨论可扩展性通常要考虑这类问题：“如果系统以某种方式增长，我们应对增长的措施有哪些”，“我们该如何添加计算资源来处理额外的负载”。

## 描述负载

首先，我们需要简洁地描述系统当前的负载，只有这样才能更好地讨论后续增长问题（例如负载加倍会意味着什么）。负载可以用称为负载参数的若干数字来描述。参数的最佳选择取决于系统的体系结构，它可能是Web服务器的每秒请求处理次数，数据库中写入的比例，聊天室的同时活动用户数量，缓存命中率等。有时平均值很重要，有时系统瓶颈来自于少数峰值。

我们以Twitter为例，使用其2012年11月发布的数据<sup>[16]</sup>。Twitter的两个典型业务操作是：

- 发布tweet消息：用户可以快速推送新消息到所有的关注者，平均大约4.6k request/sec，峰值约12k requests/sec。
- 主页时间线（Home timeline）浏览：平均300k request/sec 查看关注对象的最新消息。

仅仅处理峰值约12k的消息发送听起来并不难，但是，Twitter扩展性的挑战重点不在于消息大小，而在于巨大的扇出（fan-out）<sup>注2</sup>结构：每个用户会关注很多人，也会被很多人圈粉。此时大概有两种处理方案：

---

注2：从电子工程中借用的一个术语，它描述了输入的逻辑门连接到另一个输出门的数量。输出需要提供足够的电流来驱动所有连接的输入。在事务处理系统中，用来描述为了服务一个输入请求而需要做的请求总数。





1. 将发送的新tweet插入到全局的tweet集合中。当用户查看时间线时，首先查找所有的关注对象，列出这些人的所有tweet，最后以时间为序来排序合并。如果以图1-2的关系型数据模型，可以执行下述的SQL查询语句：

```
SELECT tweets.*, users.* FROM tweets
JOIN users ON tweets.sender_id = users.id
JOIN follows ON follows.followee_id = users.id
WHERE follows.follower_id = current_user
```

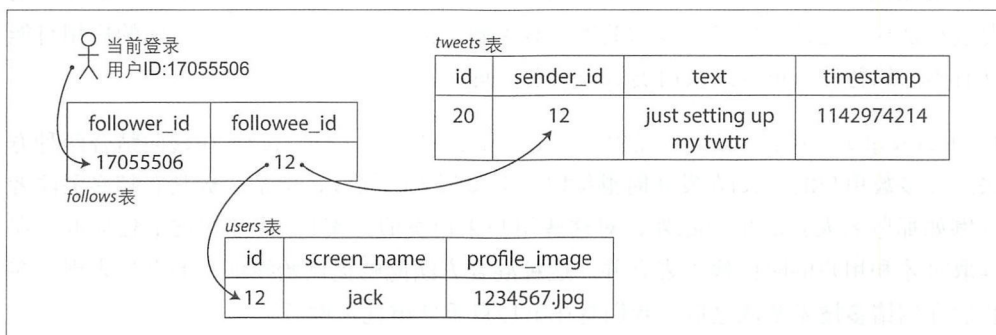


图1-2: 采用关系型数据模型来支持时间线

2. 对每个用户的时间线维护一个缓存，如图1-3所示，类似每个用户一个tweet邮箱。当用户推送新tweet时，查询其关注者，将tweet插入到每个关注者的时间线缓存中。因为已经预先将结果取出，之后访问时间线性能非常快。

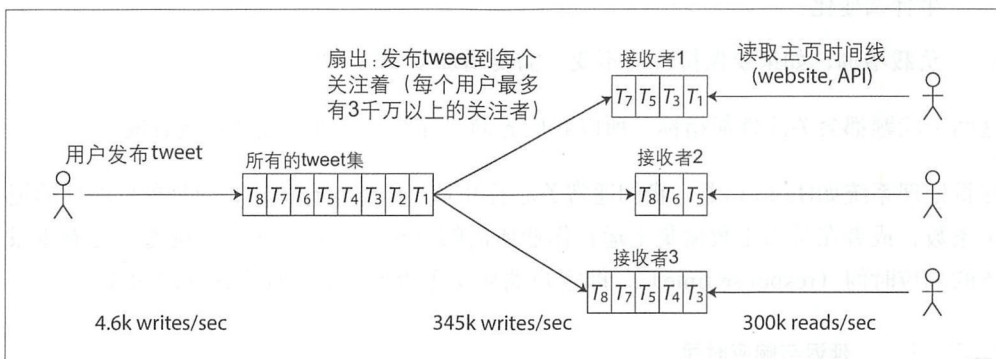


图1-3: Twitter的数据流水线方式来推送tweet，相关参数来自2012.11<sup>[16]</sup>

Twitter在其第一个版本使用了方法1，但发现主页时间线的读负载压力与日俱增，系统优化颇费周折，因此转而采用第二种方法。实践发现这样更好，因为时间线浏览tweet的压力几乎比发布tweet要高出两个数量级，基于此，在发布时多完成一些事情可以加速读性能。







然而，方法2的缺点也很明显，在发布tweet时增加了大量额外的工作。考虑平均75个关注者和每秒4.6k的tweet，则需要每秒 $4.6 \times 75 = 345\text{k}$ 速率写入缓存。但是，75这个平均关注者背后还隐藏其他事实，即关注者其实偏差巨大，例如某些用户拥有超过3000万的追随者。这就意味着峰值情况下一个tweet会导致3000万笔写入！而且要求尽量快，Twitter的设计目标是5s内完成，这成为一个巨大的挑战。

在Twitter的例子中，每个用户关注者的分布情况（还可以结合用户使用Twitter频率情况进行加权）是该案例可扩展的关键负载参数，因为它决定了扇出数。你的应用可能具有不同的特性，但可以采用类似的原则来研究具体负载。

Twitter故事最后的结局是：方法2已经得到了稳定实现，Twitter正在转向结合两种方法。大多数用户的tweet在发布时继续以一对多写入时间线，但是少数具有超多关注者（例如那些名人）的用户除外，对这些用户采用类似方案1，其推文被单独提取，在读取时才和用户的时间线主表合并。这种混合方法能够提供始终如一的良好表现。在我们介绍诸多技术基础之后，我们将在第12章重新审视该例子。

## 描述性能

描述系统负载之后，接下来设想如果负载增加将会发生什么。有两种考虑方式：

- 负载增加，但系统资源（如CPU、内存、网络带宽等）保持不变，系统性能会发生什么变化？
- 负载增加，如果要保持性能不变，需要增加多少资源？

这两个问题都会关注性能指标，所以我们先简要介绍一下如何描述系统性能。

在批处理系统如Hadoop中，我们通常关心吞吐量（throughput），即每秒可处理的记录条数，或者在某指定数据集上运行作业所需的总时间<sup>注3</sup>；而在线系统通常更看重服务的响应时间（response time），即客户端从发送请求到接收响应之间的间隔。



### 延迟与响应时间

延迟（latency）和响应时间（response time）容易混淆使用，但它们并不完全一样。通常响应时间是客户端看到的：除了处理请求时间（服务时间，service time）外，还包括来回网络延迟和各种排队延迟。延迟则是请求花费在处理上的时间。

注3：理想情况下，批量作业的运行时间是数据集的总大小除以吞吐量。在实践中，由于倾斜（数据在多个工作进程中不均匀分布）问题，系统需要等待最慢的任务完成，所以运行时间往往更长。





即使是反复发送、处理相同的请求，每次可能都会产生略微不同的响应时间。实际情况往往更复杂，由于系统要处理各种不同的请求，响应时间可能变化很大。因此，最好不要将响应时间视为一个固定的数字，而是可度量的一种数值分布。

一个例子如图1-4所示，每个灰色条表示一个服务请求，高度表示该请求的响应时间。可以看到，大多数请求是相当快的，但偶尔会有异常表示需要更长的时间。也许这些异常请求确实代价很高，例如它们的数据大很多。但有时，即使所有请求都相同，也会由于其他变量因素而引入一些随机延迟抖动，这些因素包括上下文切换和进程调度、网络数据包丢失和TCP重传、垃圾回收暂停、缺页中断和磁盘I/O，甚至服务器机架的机械振动<sup>[18]</sup>。

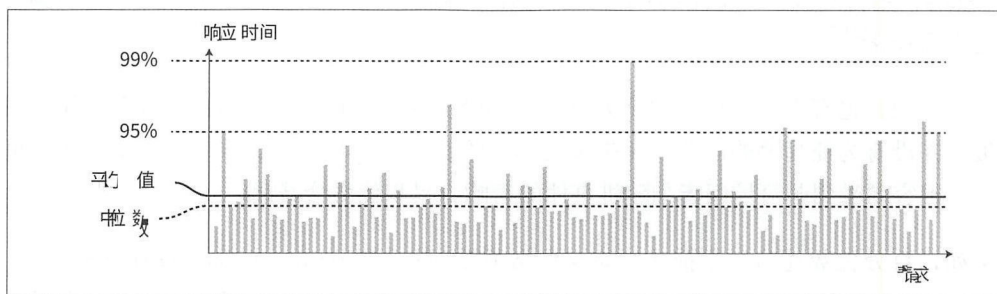


图1-4：对某服务采样100个服务请求，来说明响应时间平均值和相关百分位数

我们经常考察的是服务请求的平均响应时间（严格来说，术语“平均值”并没有明确采用何种具体公式，但通常被理解为算术平均值：给定 $n$ 个值，将所有值相加，并除以 $n$ ）。然而，如果想知道更典型的响应时间，平均值并不是合适的指标，因为它掩盖了一些信息，无法告诉有多少用户实际经历了多少延迟。

因此最好使用百分位数（percentiles）。如果已经搜集到了响应时间信息，将其从最快到最慢排序，中位数（median）就是列表中间的响应时间。例如，如果中位数响应时间为200 ms，那意味着有一半的请求响应不到200 ms，而另一半请求则需要更长的时间。

中位数指标非常适合描述多少用户需要等待多长时间：一半的用户请求的服务时间少于中位数响应时间，另一半则多于中位数的时间。因此中位数也称为50百分位数，有时缩写为p50。请注意，中位数对应单个请求；这也意味着如果某用户发了多个请求（例如包含在一个完整会话过程中，或者因为某页面包含了多个资源），那么它们中至少一个比中位数慢的概率远远大于50%。







为了弄清楚异常值有多糟糕，需要关注更大的百分位数如常见的第95、99和99.9（缩写为p95、p99和p999）值。作为典型的响应时间阈值，它们分别表示有95%、99%或99.9%的请求响应时间快于阈值。例如，如果95百分位数响应时间为1.5s，这意味着100个请求中的95个请求快于1.5s，而5个请求则需要1.5s或更长时间，如图1-4所示。

采用较高的响应时间百分位数（tail latencies，尾部延迟或长尾效应）很重要，因为它们直接影响用户的总体服务体验。例如，亚马逊采用99.9百分位数来定义其内部服务的响应时间标准，或许它仅影响1000个请求中的1个。但是考虑到请求最慢的客户往往是购买了更多的商品，因此数据量更大换言之，他们是最有价值的客户<sup>[19]</sup>。让这些客户始终保持愉悦的购物体验显然非常重要：亚马逊还注意到，响应时间每增加100ms，销售额就会下降了约1%<sup>[20]</sup>，其他研究则表明，1s的延迟增加等价于客户满意度下降16%<sup>[21,22]</sup>。

另一方面，也有人说，优化这个99.99百分位数（10 000个请求中最慢的1个）代价昂贵，并没有为亚马逊的商业目标带来足够的收益。进一步提高响应时间技术上代价更大，很容易受到非可控因素如随机事件的影响，累积优势会减弱。

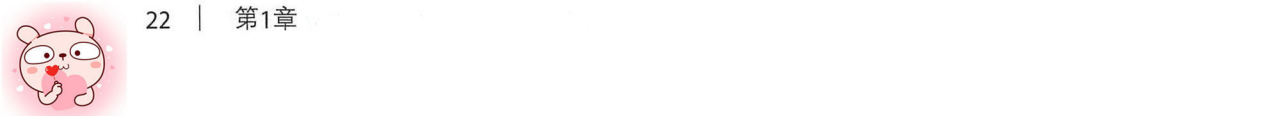
例如，百分位数通常用于描述、定义服务质量目标（Service Level Objectives，SLO）和服务质量协议（Service Level Agreements，SLA），这些是规定服务预期质量和可用性的合同。例如一份SLA合约通常会声明，响应时间中位数小于200ms，99%请求的响应时间小于1s，且要求至少99.9%的时间都要达到上述服务指标。这些指标明确了服务质量预期，并允许客户在不符合SLA的情况下进行赔偿。

排队延迟往往在高百分数响应时间中影响很大。由于服务器并行处理的请求有限（例如，CPU内核数的限制），正在处理的少数请求可能会阻挡后续请求，这种情况有时被称为队头阻塞。即使后续请求可能处理很简单，但它阻塞在等待先前请求的完成，客户端将会观察到极慢的响应时间。因此，很重要的一点是要在客户端来测量响应时间。

因此，当测试系统可扩展性而人为产生负载时，负载生成端要独立于响应时间来持续发送请求。如果客户端在发送请求之前总是等待先前请求的完成，就会在测试中人为地缩短了服务器端的累计队列深度，这就带来了测试偏差<sup>[23]</sup>。

## 应对负载增加的方法

我们已经讨论了描述负载的参数以及衡量性能的相关指标，接下来讨论可扩展性：即当负载参数增加时，应如何保持良好性能？





## 实践中的百分位数

对于后台服务，如果一次完整的服务里包含了多次请求调用，此时高百分位数指标尤为重要。即使这些子请求是并行发送、处理，但最终用户仍然需要等待最慢的那个调用完成才行。如图1-5所示，哪怕一个缓慢的请求处理，即可拖累整个服务。即使只有很小百分比的请求缓慢，如果某用户总是频繁产生这种调用，最终总体变慢的概率就会增加（即长尾效应<sup>[24]</sup>）。

最好将响应时间百分位数添加到服务系统监控中，持续跟踪该指标。例如，设置一个10min的滑动窗口，监控其中响应时间，滚动计算窗口中的中位数和各种百分位数，然后绘制性能图表。

一种简单的实现方案是在时间窗口内保留所有请求的响应时间列表，每分钟做一次排序。如果这种方式效率太低，可以采用一些近似算法（如正向衰减<sup>[25]</sup>，t-digest<sup>[26]</sup>或HdrHistogram<sup>[27]</sup>）来计算百分位数，其CPU和内存开销很低。同时请注意，降低采样时间精度或直接组合来自多台机器的数据，在数学上没有太大意义，聚合响应时间的正确方法是采用直方图<sup>[28]</sup>。

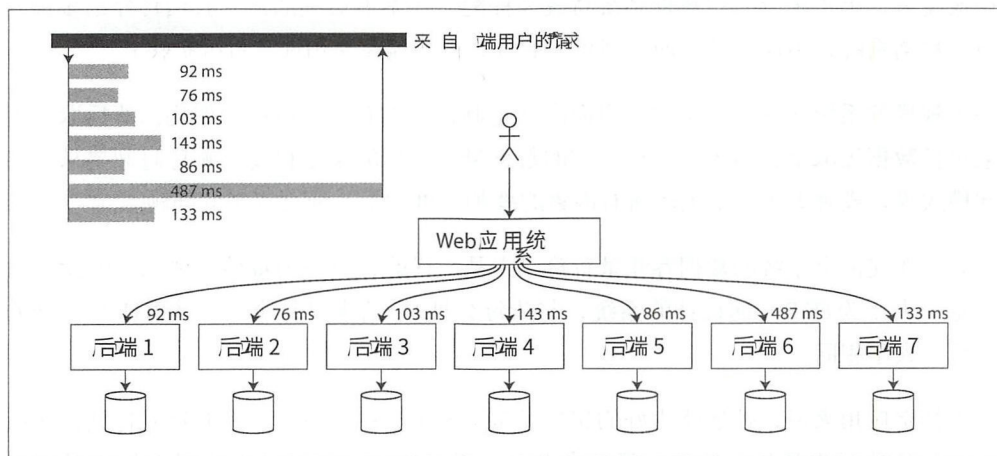


图1-5：一个服务涉及多个不同的后端调用，则最慢的调用会拖累整个服务响应时间

首先，针对特定级别负载而设计的架构不太可能应付超出预设目标10倍的实际负载。如果目标服务处于快速增长阶段，那么需要认真考虑每增加一个数量级的负载，架构应如何设计。







现在谈论更多的是如何在垂直扩展（即升级到更强大的机器）和水平扩展（即将负载分布到多个更小的机器）之间做取舍。在多台机器上分配负载也被称为无共享体系结构。在单台机器上运行的系统通常更简单，然而高端机器可能非常昂贵，且扩展水平有限，最终往往还是无法避免需要水平扩展。实际上，好的架构通常要做些实际取舍，例如，使用几个强悍的服务器仍可以比大量的小型虚拟机来得更简单、便宜。

某些系统具有弹性特征，它可以自动检测负载增加，然后自动添加更多计算资源，而其他系统则是手动扩展（人工分析性能表现，之后决定添加更多计算）。如果负载高度不可预测，则自动弹性系统会更加高效，但或许手动方式可以减少执行期间的意外情况（参阅第6章的“分区再平衡”）。

把无状态服务分布然后扩展至多台机器相对比较容易，而有状态服务从单个节点扩展到分布式多机环境的复杂性会大大增加。出于这个原因，直到最近通常的做法一直是，将数据库运行在一个节点上（采用垂直扩展策略），直到高扩展性或高可用性的要求迫使不得不做水平扩展。

然而，随着相关分布式系统专门组件和编程接口越来越好，至少对于某些应用类型来讲，上述通常做法或许会发生改变。可以乐观设想，即使应用可能并不会处理大量数据或流量，但未来分布式数据系统将成为标配。在本书后续部分，我们将介绍多种分布式数据系统，不仅可以帮助提高可扩展性，也会提高易用性与可维护性。

超大规模的系统往往针对特定应用而高度定制，很难有一种通用的架构。背后取舍因素包括数据读取量、写入量、待存储的数据量、数据的复杂程度、响应时间要求、访问模式等，或者更多的是上述所有因素的叠加，再加上其他更复杂的问题。

例如，即使两个系统的数据吞吐量折算下来是一样的，但是为每秒处理100000次请求（每个大小为1KB）而设计的系统，与为每分钟3个请求（每个大小为2GB）设计的系统会大不相同。

对于特定应用来说，扩展能力好的架构通常会做出某些假设，然后有针对性地优化设计，如哪些操作是最频繁的，哪些负载是少数情况。如果这些假设最终发现是错误的，那么可扩展性的努力就白费了，甚至会出现与设计预期完全相反的情况。对于早期的初创公司或者尚未定型的产品，快速迭代推出产品功能往往比投入精力来应对不可知的扩展性更为重要。

可扩展架构通常都是从通用模块逐步构建而来，背后往往有规律可循，所以本书将会讨论这些通用模块和常见模式，希望对读者有所借鉴。





## 可维护性

众所周知，软件的大部分成本并不在最初的开发阶段，而是在于整个生命周期内持续的投入，这包括维护与缺陷修复，监控系统来保持正常运行、故障排查、适配新平台、搭配新场景、技术缺陷的完善以及增加新功能等。

不幸的是，许多从业者根本不喜欢维护这些所谓的遗留系统，例如修复他人埋下的错误，或者使用过时的开发平台，或者被迫做不喜欢的工作。坦白说，每一个遗留系统总有其过期的理由，所以很难给出一个通用的建议该如何处理它们。

但是，换个角度，我们可以从软件设计时开始考虑，尽可能较少维护期间的麻烦，甚至避免造出容易过期的系统。为此，我们将特别关注软件系统的三个设计原则：

### 可运维性

方便运营团队来保持系统平稳运行。

### 简单性

简化系统复杂性，使新工程师能够轻松理解系统。注意这与用户界面的简单性并不一样。

### 可演化性

后续工程师能够轻松地对系统进行改进，并根据需求变化将其适配到非典型场景，也称为可延伸性、易修改性或可塑性。

与可靠性和可扩展性类似，实现上述这些目标也没有简单的解决方案。接下来，我们首先建立对这三个特性的理解。

## 可运维性：运维更轻松

有人认为，“良好的操作性经常可以化解软件的局限性，而不规范的操作则可以轻松击垮软件<sup>[12]</sup>”。虽然某些操作可以而且应该是自动化的，但最终还是需要人来执行配置并确保正常工作。

运营团队对于保持软件系统顺利运行至关重要。一个优秀的运营团队通常至少负责以下内容<sup>[29]</sup>：

- 监视系统的健康状况，并在服务出现异常状态时快速恢复服务。
- 追踪问题的原因，例如系统故障或性能下降。







- 保持软件 and 平台至最新状态，例如安全补丁方面。
- 了解不同系统如何相互影响，避免执行带有破坏性的操作。
- 预测未来可能的问题，并在问题发生之前即使解决（例如容量规划）。
- 建立用于部署、配置管理等良好的实践规范和工具包。
- 执行复杂的维护任务，例如将应用程序从一个平台迁移到另一个平台。
- 当配置更改时，维护系统的安全稳健。
- 制定流程来规范操作行为，并保持生产环境稳定。
- 保持相关知识的传承（如对系统理解），例如发生团队人员离职或者新员工加入等。

良好的可操作性意味着使日常工作变得简单，使运营团队能够专注于高附加值的任务。数据系统设计可以在这方面贡献很多，包括：

- 提供对系统运行时行为和内部的可观测性，方便监控。
- 支持自动化，与标准工具集成。
- 避免绑定特定的机器，这样在整个系统不间断运行的同时，允许机器停机维护。
- 提供良好的文档和易于理解的操作模式，诸如“如果我做了X，会发生Y”。
- 提供良好的默认配置，且允许管理员在需要时方便地修改默认值。
- 尝试自我修复，在需要时让管理员手动控制系统状态。
- 行为可预测，减少意外发生。

## 简单性：简化复杂度

小型软件项目通常可以写出简单而漂亮的代码，但随着项目越来越大，就会越来越复杂和难以理解。这种复杂性拖慢了开发效率，增加了维护成本。一个过于复杂的软件项目有时被称为一个“大泥潭”<sup>[30]</sup>。

复杂性有各种各样的表现方式：状态空间的膨胀，模块紧耦合，令人纠结的相互依赖关系，不一致的命名和术语，为了性能而采取的特殊处理，为解决某特定问题而引入的特殊框架等。在参考文献<sup>[31-33]</sup>中有很多这方面的讨论。

复杂性使得维护变得越来越困难，最终会导致预算超支和开发进度滞后。对于复杂的软件系统，变更而引入潜在错误的风险会显著加大，最终开发人员更加难以准确理





解、评估或者更加容易忽略相关系统行为，包括背后的假设，潜在的后果，设计之外的模块交互等。相反，降低复杂性可以大大提高软件的可维护性，因此简单性应该是我们构建系统的关键目标之一。

简化系统设计并不意味着减少系统功能，而主要意味着消除意外方面的复杂性，正如Moseley和Marks<sup>[32]</sup>把复杂性定义为一种“意外”，即它并非软件固有、被用户所见或感知，而是实现本身所衍生出来的问题。

消除意外复杂性最好手段之一是抽象。一个好的设计抽象可以隐藏大量的实现细节，并对外提供干净、易懂的接口。一个好的设计抽象可用于各种不同的应用程序。这样，复用远比多次重复实现更有效率；另一方面，也带来更高质量的软件，而质量过硬的抽象组件所带来的好处，可以使运行其上的所有应用轻松获益。

例如，高级编程语言作为一种抽象，可以隐藏机器汇编代码、CPU寄存器和系统调用等细节和复杂性。SQL作为一种抽象，隐藏了内部复杂的磁盘和内存数据结构，以及来自多客户端的并发请求，系统崩溃之后的不一致等问题。当然，使用高级编程语言最终并没有脱离机器汇编代码，只是并非直接使用，与汇编代码打交道的事情已经由编程语言抽象为高效接口代我们完成。

然而，设计好的抽象还是很有挑战性。在分布式系统领域中，虽然已有许多好的算法可供参考，但很多时候我们并不太清楚究竟该如何利用他们，封装到抽象接口之中，最终帮助将系统的复杂性降低到可掌控的级别。

本书我们将广泛考察如何设计好的抽象，这样至少能够将大型系统的一部分抽象为定义明确、可重用的组件。

## 可演化性：易于改变

一成不变的系统需求几乎没有，想法和目标经常在不断变化：适配新的外部环境，新的用例，业务优先级的变化，用户要求的新功能，新平台取代旧平台，法律或监管要求的变化，业务增长促使架构的演变等。

在组织流程方面，敏捷开发模式为适应变化提供了很好的参考。敏捷社区还发布了很多技术工具和模式，以帮助在频繁变化的环境中开发软件，例如测试驱动开发（TDD）和重构。

这些敏捷开发技术目前多数还只是针对小规模、本地模式（例如同一应用程序中的几个源代码文件）环境。本书将探索在更大的数据系统层面上提高敏捷性，系统由多个







不同特性的应用或者服务协作而成。例如，对于Twitter的案例（参见本章前面的“描述负载”），如何从方法1过渡到方法2，重构Twitter架构来实现主页时间线。

我们的目标是可以轻松地修改数据系统，使其适应不断变化的需求，这和简单性与抽象性密切相关：简单易懂的系统往往比复杂的系统更容易修改。这是一个非常重要的理念，我们将采用另一个不同的词来指代数据系统级的敏捷性，即可演化性<sup>[34]</sup>。

## 小结

这一章我们探讨了一些关于数据密集型应用的基本原则，这些原则将指导如何阅读本书的其余部分。

一个应用必须完成预期的多种需求，主要包括功能性需求（即应该做什么，比如各种存储、检索、搜索和处理数据）和一些非功能性需求（即常规特性、例如安全性、可靠性、合规性、可伸缩性、兼容性和可维护性）。本章我们着重梳理讨论了可靠性、可扩展性和可维护性。

可靠性意味着即使发生故障，系统也可以正常工作。故障包括硬件（通常是随机的，不相关的）、软件（缺陷通常是系统的，更加难以处理）以及人为（总是很难避免时不时会出错）方面。容错技术可以很好地隐藏某种类型故障，避免影响最终用户。

可扩展性是指负载增加时，有效保持系统性能的相关技术策略。为了讨论可扩展性，我们首先探讨了如何定量描述负载和性能。简单地以Twitter浏览时间线为例描述负载，并将响应时间百分位数作为衡量性能的有效方式。对于可扩展的系统，增加处理能力的同时，还可以在高负载情况下持续保持系统的高可靠性。

可维护性则意味着许多方面，但究其本质是为了让工程和运营团队更为轻松。良好的抽象可以帮助降低复杂性，并使系统更易于修改和适配新场景。良好的可操作性意味着对系统健康状况有良好的可观测性和有效的管理方法。

然而知易行难，使应用程序可靠、可扩展或可维护并不容易。考虑到一些重要的模式和技术在很多不同应用中普遍适用，在接下来的几章中，我们就一些数据密集系统例子，分析它们如何实现上述这些目标。

在本书的第三部分，我们将看到更多如图1-1所示的包含多个组件但更为复杂的例子。





## 参考文献

- [1] Michael Stonebraker and Uğur Çetintemel: “ ‘One Size Fits All’ : An Idea Whose Time Has Come and Gone,” at *21st International Conference on Data Engineering (ICDE)*, April 2005.
- [2] Walter L. Heimerdinger and Charles B. Weinstock: “A Conceptual Framework for System Fault Tolerance,” Technical Report CMU/SEI-92-TR-033, Software Engineering Institute, Carnegie Mellon University, October 1992.
- [3] Ding Yuan, Yu Luo, Xin Zhuang, et al.: “Simple Testing Can Prevent Most Critical Failures: An Analysis of Production Failures in Distributed Data-Intensive Systems,” at *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, October 2014.
- [4] Yury Izrailevsky and Ariel Tseitlin: “The Netflix Simian Army,” *techblog.netflix.com*, July 19, 2011.
- [5] Daniel Ford, François Labelle, Florentina I. Popovici, et al.: “Availability in Globally Distributed Storage Systems,” at *9th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, October 2010.
- [6] Brian Beach: “Hard Drive Reliability Update-Sep 2014,” *backblaze.com*, September 23, 2014.
- [7] Laurie Voss: “AWS: The Good, the Bad and the Ugly,” *blog.awe.sm*, December 18, 2012.
- [8] Haryadi S. Gunawi, Mingzhe Hao, Tanakorn Leesatapornwongsa, et al.: “What Bugs Live in the Cloud?,” at *5th ACM Symposium on Cloud Computing (SoCC)*, November 2014. doi:10.1145/2670979.2670986.
- [9] Nelson Minar: “Leap Second Crashes Half the Internet,” *somebits.com*, July 3, 2012.
- [10] Amazon Web Services: “Summary of the Amazon EC2 and Amazon RDS Service Disruption in the US East Region,” *aws.amazon.com*, April 29, 2011.







- [11] Richard I. Cook: “How Complex Systems Fail,” Cognitive Technologies Laboratory, April 2000.
- [12] Jay Kreps: “Getting Real About Distributed System Reliability,” *blog.empathybox.com*, March 19, 2012.
- [13] David Oppenheimer, Archana Ganapathi, and David A. Patterson: “Why Do Internet Services Fail, and What Can Be Done About It?,” at *4th USENIX Symposium on Internet Technologies and Systems* (USITS), March 2003.
- [14] Nathan Marz: “Principles of Software Engineering, Part 1,” *nathanmarz.com*, April 2, 2013.
- [15] Michael Jurewitz: “The Human Impact of Bugs,” *jury.me*, March 15, 2013.
- [16] Raffi Krikorian: “Timelines at Scale,” at *QCon San Francisco*, November 2012.
- [17] Martin Fowler: *Patterns of Enterprise Application Architecture*. Addison Wesley, 2002. ISBN: 978-0-321-12742-6.
- [18] Kelly Sommers: “After all that run around, what caused 500ms disk latency even when we replaced physical server?” *twitter.com*, November 13, 2014.
- [19] Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, et al.: “Dynamo: Amazon’s Highly Available Key-Value Store,” at *21st ACM Symposium on Operating Systems Principles* (SOSP), October 2007.
- [20] Greg Linden: “Make Data Useful,” slides from presentation at Stanford University Data Mining class (CS345), December 2006.
- [21] Tammy Everts: “The Real Cost of Slow Time vs Downtime,” *webperformance-today.com*, November 12, 2014.
- [22] Jake Brutlag: “Speed Matters for Google Web Search,” *googleresearch.blogspot.co.uk*, June 22, 2009.
- [23] Tyler Treat: “Everything You Know About Latency Is Wrong,” *bravenewgeek.com*, December 12, 2015.
- [24] Jeffrey Dean and Luiz André Barroso: “The Tail at Scale,” *Communications of the*





ACM, volume 56, number 2, pages 74-80, February 2013. doi: 10.1145/2408776.2408794.

[25] Graham Cormode, Vladislav Shkapenyuk, Divesh Srivastava, and Bojian Xu: “Forward Decay: A Practical Time Decay Model for Streaming Systems,” at *25th IEEE International Conference on Data Engineering (ICDE)*, March 2009.

[26] Ted Dunning and Otmar Ertl: “Computing Extremely Accurate Quantiles Using t-Digests,” *github.com*, March 2014.

[27] Gil Tene: “HdrHistogram,” *hdrhistogram.org*.

[28] Baron Schwartz: “Why Percentiles Don’t Work the Way You Think,” *vividcortex.com*, December 7, 2015.

[29] James Hamilton: “On Designing and Deploying Internet-Scale Services,” at *21st Large Installation System Administration Conference (LISA)*, November 2007.

[30] Brian Foote and Joseph Yoder: “Big Ball of Mud,” at *4th Conference on Pattern Languages of Programs (PLoP)*, September 1997.

[31] Frederick P Brooks: “No Silver Bullet – Essence and Accident in Software Engineering,” in *The Mythical Man-Month*, Anniversary edition, Addison-Wesley, 1995. ISBN: 978-0-201-83595-3.

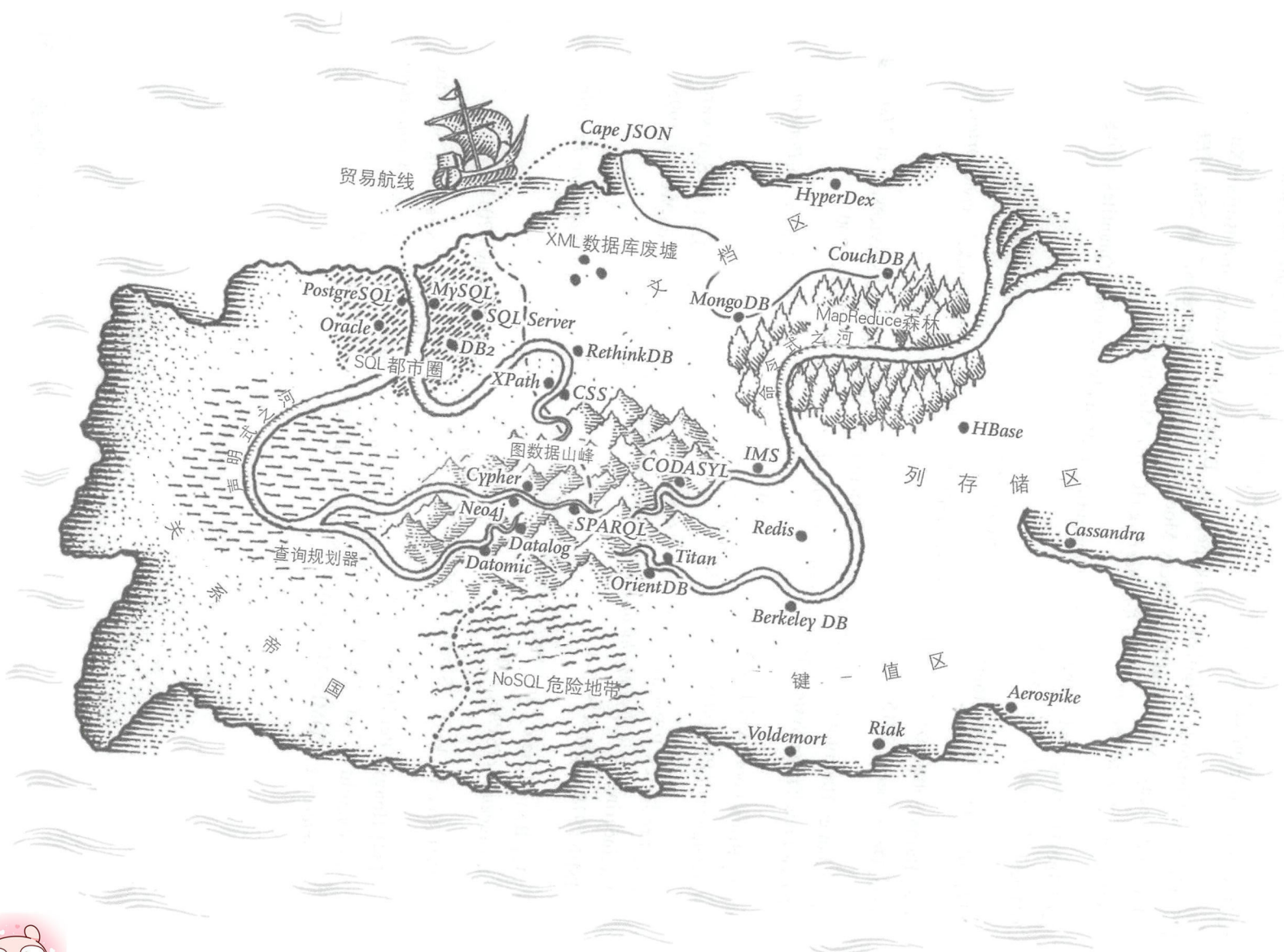
[32] Ben Moseley and Peter Marks: “Out of the Tar Pit,” at *BCS Software Practice Advancement (SPA)*, 2006.

[33] Rich Hickey: “Simple Made Easy,” at *Strange Loop*, September 2011.

[34] Hongyu Pei Breivold, Ivica Crnkovic, and Peter J. Eriksson: “Analyzing Software Evolvability,” at *32nd Annual IEEE International Computer Software and Applications Conference (COMPSAC)*, July 2008. doi:10.1109/COMPSAC.2008.50.









# 数据模型与查询语言

语言的边界就是世界的边界。

——Ludwig Wittgenstein, 《逻辑哲学论》(1922)

数据模型可能是开发软件最重要的部分，它们不仅对软件的编写方式，而且还对如何思考待解决的问题都有深远的影响。

大多数应用程序是通过一层一层叠加数据模型来构建的。每一层都面临的关键问题是：如何将其用下一层来表示？例如：

1. 作为一名应用程序开发人员，观测现实世界（其中包括人员、组织、货物、行为、资金流动、传感器等），通过对象或数据结构，以及操作这些数据结构的API来对其建模。这些数据结构往往特定于该应用。
2. 当需要存储这些数据结构时，可以采用通用数据模型（例如JSON或XML文档、关系数据库中的表或图模型）来表示。
3. 数据库工程师接着决定用何种内存、磁盘或网络的字节格式来表示上述JSON/XML/关系/图形数据。数据表示需要支持多种方式的查询、搜索、操作和处理数据。
4. 在更下一层，硬件工程师则需要考虑用电流、光脉冲、磁场等来表示字节。

复杂的应用程序可能会有更多的中间层，例如基于API来构建上层API，但是基本思想相同：每层都通过提供一个简洁的数据模型来隐藏下层的复杂性。这些抽象机制使得不同的人群可以高效协作，例如数据厂商的工程师和使用数据库的应用程序开发人员一起合作。







有许多不同类型的数据模型，每种数据模型都有其最佳使用的若干假设。有些用法很简单，有些则根本不支持；有些操作性能很快，有些则很差；有些数据转换非常自然，有些则异常笨拙。

精通一种数据模型都需要很大功夫（试想有多少关于关系数据建模的书籍）。即使只使用一种数据模型且不用担心内部机制，构建软件也很有挑战。然而考虑到数据模型对其上的软件应用有着巨大影响（哪些可以做和哪些不能做），因此需要慎重选择合适的数据库模型。

本章将介绍一系列用于数据存储和查询的通用数据库模型（上一列表中的第2点）。特别地，我们将比较关系模型、文档模型和一些基于图的数据模型。我们还将讨论多种查询语言并比较它们的使用场景。第3章将讨论存储引擎，即这些数据库模型是如何实现的（列表的第3点）。

## 关系模型与文档模型

现在最著名的数据库模型可能是SQL，它基于Edgar Codd于1970年提出的关系模型<sup>[1]</sup>：数据被组织成关系（relations），在SQL中称为表（table），其中每个关系都是元组（tuples）的无序集合（在SQL中称为行）。

关系模型曾经只是一个理论建议，当时很多人怀疑它是否能够被高效地实现。然而，到了20世纪80年代中期，关系数据库管理系统（RDBMS）和SQL已经成为大多数需要存储、查询具有某种规则结构数据的首选工具。关系数据库的主导地位已经持续了25~30年，这在计算机历史称得上一段不朽传奇。

关系数据库的核心在于商业数据处理，20世纪60年代和70年代主要运行在大型计算机之上。从今天的角度来看，用例看起来很常见，主要是事务处理（包括输入销售和银行交易、航空公司订票、仓库库存）和批处理（例如客户发票、工资单、报告）。

当时的其他数据库迫使应用开发人员考虑数据的内部表示。关系模型的目标就是将实现细节隐藏在更简洁的接口后面。

多年来，在数据存储和查询方面存在着其他许多竞争技术。在20世纪70年代和80年代初期，网络模型和层次模型是两个主要的选择，但最终关系模型主宰了这个领域。对象数据库曾在20世纪80年代后期和90年代初期起起伏伏。XML数据库则出现在21世纪初，但也仅限于利基市场。关系模型的每个竞争者都曾聒噪一时，可惜无一持久<sup>[2]</sup>。





随着计算机变得越来越强大和网络化，服务目的日益多样化。值得注意的是，关系数据库超出了它们最初的商业数据处理范围，顺利推广到了各种各样的用例。当前在网上看到的大部分内容很多仍然是由关系数据库所支撑的，无论是在线发布、论坛、社交网络、电子商务、游戏、SaaS等。

## NoSQL的诞生

进入21世纪，NoSQL成为推翻关系模式主导地位的又一个竞争者。“NoSQL”这个名字是不恰当的，因为它其实并不代表具体的某些技术，它最初只是作为一个吸引人眼球的Twitter标签频频出现在2009年的开源、分布式以及非关系数据库的见面会上<sup>[3]</sup>。尽管如此，这个称呼还是让人有所触动，并迅速传遍了网络创业社区。现在很多新兴的数据库系统总是会打上NoSQL的标签，而其含义也已经被逆向重新解释为“不仅仅是SQL”<sup>[4]</sup>。

采用NoSQL数据库有这样几个驱动因素，包括：

- 比关系数据库更好的扩展性需求，包括支持超大数据集或超高写入吞吐量。
- 普遍偏爱免费和开源软件而不是商业数据库产品。
- 关系模型不能很好地支持一些特定的查询操作。
- 对关系模式一些限制性感到沮丧，渴望更具动态和表达力的数据模型<sup>[5]</sup>。

不同的应用程序有不同的需求，某个用例的最佳的技术选择未必适合另一个用例。因此，在可预见的将来，关系数据库可能仍将继续与各种非关系数据存储一起使用，这种思路有时也被称为混合持久化<sup>[3]</sup>。

## 对象-关系不匹配

现在大多数应用开发都采用面向对象的编程语言，由于兼容性问题，普遍对SQL数据模型存在抱怨：如果数据存储的关系表中，那么应用层代码中的对象与表、行和列的数据库模型之间需要一个笨拙的转换层。模型之间的脱离有时被称为阻抗失谐<sup>注1</sup>。

ActiveRecord和Hibernate这样的对象-关系映射（ORM）框架则减少了此转换层所需的样板代码量，但是他们并不能完全隐藏两个模型之间的差异。

---

注1：从电子学借用的一个术语。每个电路的输入和输出都有一定的阻抗（交流电阻）。将一个电路的输出连接到另一个电路的输入时，如果两个电路的输出和输入阻抗匹配，则连接上的功率传输将被最大化。阻抗不匹配会导致信号反射和其他故障。







例如，图2-1展示了如何在关系模式中表示简历（类似LinkedIn profile）。整个简历可以通过唯一的标识符`user_id`来标识。像`first_name`和`last_name`这样的字段在每个用户中只出现一次，所以可以将其建模为`users`表中的列。然而，大多数人在他们的职业（职位）中有一个以上的工作，并且可能有多个教育阶段和任意数量的联系信息。用户与这些项目之间存在一对多的关系，可以用多种方式来表示：

- 在传统的SQL模型（SQL：1999之前）中，最常见的规范化表示是将职位、教育和联系信息放在单独的表中，并使用外键引用`users`表，如图2-1所示。
- 之后的SQL标准增加了对结构化数据类型和XML数据的支持。这允许将多值数据存储在单行内，并支持在这些文档中查询和索引。Oracle、IBM DB2、MS SQL Server和PostgreSQL都不同程度上支持这些功能<sup>[6,7]</sup>。一些数据库也支持JSON数据类型，例如IBM DB2、MySQL和PostgreSQL<sup>[8]</sup>。
- 第三个选项是将工作、教育和联系信息编码为JSON或XML文档，将其存储在数据库的文本列中，并由应用程序解释其结构和内容。对于此方法，通常不能使用数据库查询该编码列中的值。

对于像简历这样的数据结构，它主要是一个自包含的文档（document），因此用JSON表示非常合适，参见示例2-1。与XML相比，JSON的吸引力在于它更简单。面向文档的数据库（如MongoDB<sup>[9]</sup>、RethinkDB<sup>[10]</sup>、CouchDB<sup>[11]</sup>和Espresso<sup>[12]</sup>）都支持该数据模型。

示例2-1：将LinkedIn简历表示为JSON文档

```
{
  "user_id":      251,
  "first_name":   "Bill",
  "last_name":    "Gates",
  "summary":      "Co-chair of the Bill & Melinda Gates... Active blogger.",
  "region_id":    "us:91",
  "industry_id":  131,
  "photo_url":    "/p/7/000/253/05b/308dd6e.jpg",
  "positions": [
    {"job_title": "Co-chair", "organization": "Bill & Melinda Gates Foundation"},
    {"job_title": "Co-founder, Chairman", "organization": "Microsoft"}
  ],
  "education": [
    {"school_name": "Harvard University", "start": 1973, "end": 1975},
    {"school_name": "Lakeside School, Seattle", "start": null, "end": null}
  ],
  "contact_info": {
    "blog": "http://thegatesnotes.com",
    "twitter": "http://twitter.com/BillGates"
  }
}
```

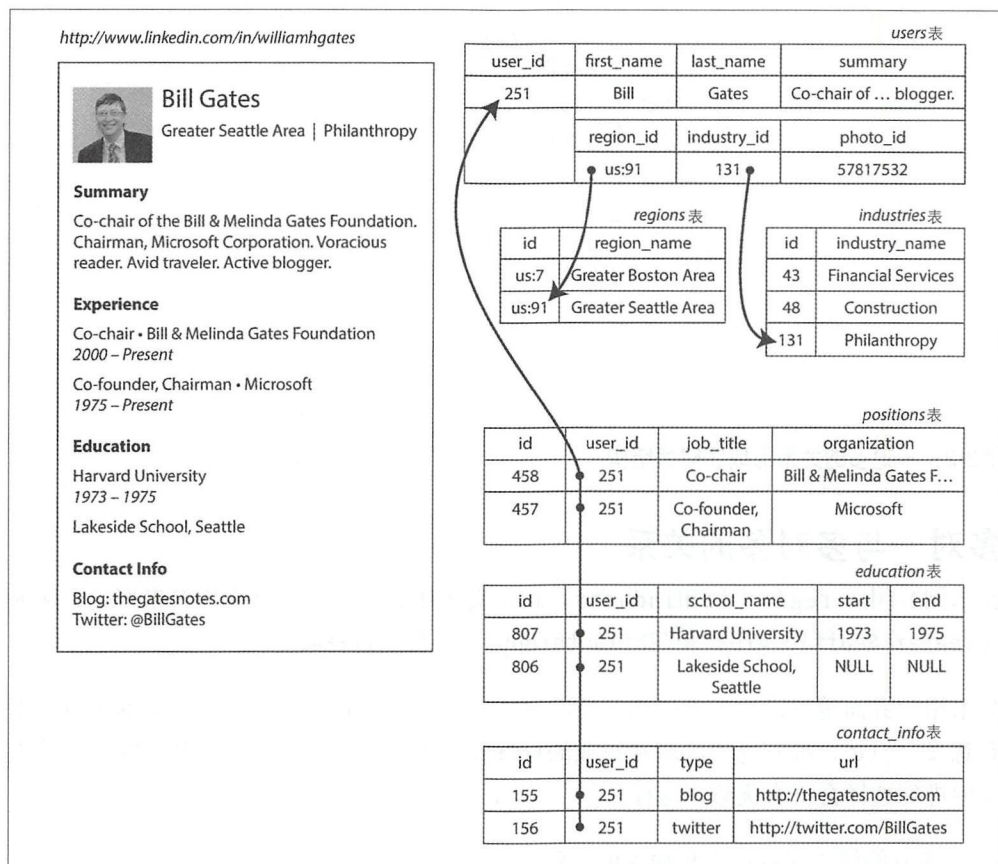


图2-1：使用关系模式来表示LinkedIn简历（比尔·盖茨的照片由维基共享提供）

一些开发人员认为JSON模型减少了应用程序代码和存储层之间的阻抗失配。然而，正如我们将在第4章中看到的那样，JSON作为数据编码格式也存在问题。缺乏模式常常被认为是一个优势；我们将在本章后面“文档模型的模式灵活性”中讨论这个问题。

JSON表示比图2-1的多表模式具有更好的局部性。如果要在关系模式中读取一份简历，那么要么执行多个查询（通过user\_id查询每个表），要么在users表及其从属表之间执行混乱的多路联结。而对于JSON表示方法，所有的相关信息都在一个地方，一次查询就够了。

用户简历到用户的职位、教育历史和联系信息的一对多关系，意味着数据存在树状结构，JSON表示将该树结构显式化（见图2-2）。



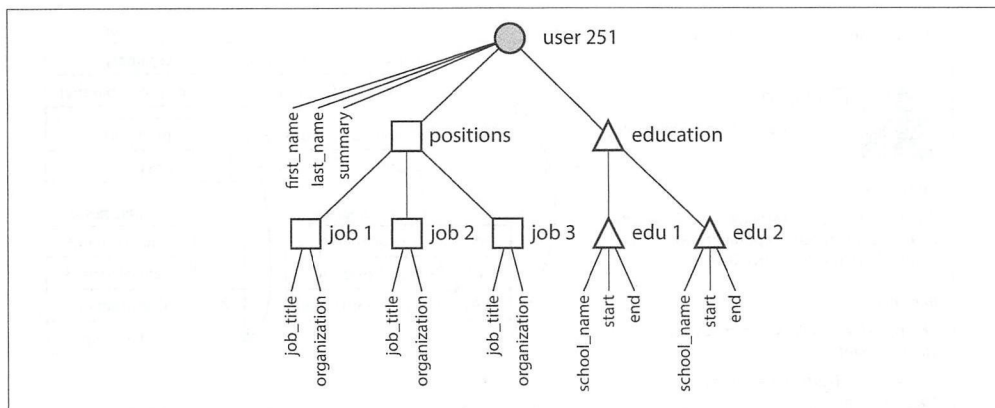


图2-2：一对多的关系形成了树状结构

## 多对一与多对多的关系

在示例2-1中，`region_id`和`industry_id`定义为ID，而不是纯文本字符串形式，例如“Greater Seattle Area”和“Philanthropy”。为什么这样做呢？

如果用户界面是可以输入地区或行业的自由文本字段，则将其存储为纯文本字符串更有意义。但是，拥有地理区域和行业的标准化列表，并让用户从下拉列表或自动填充器中进行选择会更有优势，这样：

- 所有的简历保持样式和输入值一致。
- 避免歧义（例如，如果存在一些同名的城市）。
- 易于更新：名字只保存一次，因此，如果需要改变（例如，由于政治事件而更改城市名称），可以很容易全面更新。
- 本地化支持：当网站被翻译成其他语言时，标准化的列表可以方便本地化，因此地区 and 行业可以用查看者的母语来显示。
- 更好的搜索支持：例如，搜索华盛顿州的慈善家可以匹配到这个简历，这是因为地区列表可以将西雅图属于华盛顿的信息编码进来（而从“大西雅图地区”字符串中并不能看出来西雅图属于华盛顿）。

无论是存储ID还是文本字符串，都涉及内容重复的问题。当使用ID时，对人类有意义的信息（例如慈善这个词）只存储在一个地方，引用它的所有内容都使用ID（ID只在数据库中有意义）。当直接存储文本时，则使用它的每条记录中都保存了一份这样可读信息。



使用ID的好处是，因为它对人类没有任何直接意义，所以永远不需要直接改变：即使ID标识的信息发生了变化，它也可以保持不变。任何对人类有意义的东西都可能在将来某个时刻发生变更。如果这些信息被复制，那么所有的冗余副本也都需要更新。这会导致更多写入开销，并且存在数据不一致的风险（信息的一些副本被更新，而其他副本未更新）。消除这种重复正是数据库规范化的核心思想<sup>注2</sup>。



数据库管理员和开发人员有时喜欢争论规范化与反规范化，此处我们暂不做评论。在本书的第三部分，我们将重回这个话题，并探索处理缓存、反规范化和派生数据更为系统化的方法。

然而这种数据规范化需要表达多对一的关系（许多人生活在同一地区，许多人在同一行业工作），这并不是很适合文档模型。对于关系数据库，由于支持联结操作，可以很方便地通过ID来引用其他表中的行。而在文档数据库中，一对多的树状结构不需要联结，支持联结通常也很弱<sup>注3</sup>。

如果数据库本身不支持联结，则必须在应用程序代码中，通过对数据库进行多次查询来模拟联结（对于上述例子，地区和行业的列表很小且一段时间内不太可能发生变化，应用程序可以简单地将它们缓存在内存中。但无论如何，联结的工作其实从数据库转移到了应用层）。

即使应用程序的初始版本非常适合采用无联结的文档模型，但随着应用支持越来越多的功能，数据也变得更加互联一体化。例如，考虑以下我们可能对简历进行的修改或扩充：

### 组织和学校作为实体

在前面的定义中，`organization`（用户所在的公司）和`school_name`（用户所在的学校）都是字符串。也许他们应该定义为实体的引用？然后每个组织、学校或大学都可以拥有自己的网页（logo、新闻发布源等）。每个简历都可以链接到相关组织和学校，包括他们的logo和其他信息（来自LinkedIn的一个例子，见图2-3）。

---

注2：关于关系模型的文献还区分了多种不同的范式，但这些区分的实际价值上并不大。一个经验法则是，如果复制了多份重复的数据，那么该模式通常就违背了规范化。

注3：在编写本书时，RethinkDB支持联结，MongoDB不支持联结，在CouchDB中，只有预先声明的视图支持联结。





## 推荐

假设添加这样一个新功能：一个用户可以推荐其他用户。推荐显示在被推荐者的简历上，并附上推荐人的姓名和照片。如果推荐人更新了他们的照片，他们所写的任何推荐都需要显示新照片。因此，推荐需要有一个到作者简历的引用。

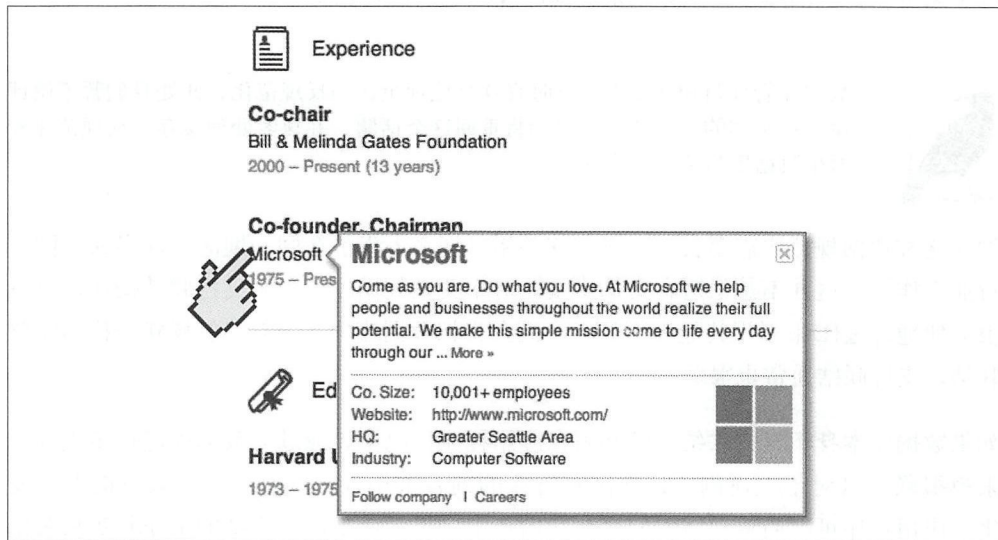


图2-3：公司名称不是一个字符串，而是一个指向公司实体的链接（上述例子为linkedin.com截图）

图2-4展示了这些新功能如何定义多对多的关系。每个虚线矩形框内的数据可以组织为一个文档，但是指向组织、学校以及其他用户的关系则需要表示为引用，并且在查询时需要联结操作。

## 文档数据库是否在重演历史？

虽然关系数据库中经常使用多对多的关系和联结，但文档数据库和NoSQL再次引发了如何最佳表示数据关系的争论。而类似争论并不是第一次，事实上，它可以追溯到最早的计算机数据库系统。

20世纪70年代最受欢迎的商业数据处理数据库是IBM信息管理系统（Information Management System, IMS），最初是为了阿波罗太空计划中的库存管理而开发的，并于1968年首次商业发布<sup>[13]</sup>。至今它仍在维护和使用，通常运行在IBM大型机OS / 390<sup>[14]</sup>。

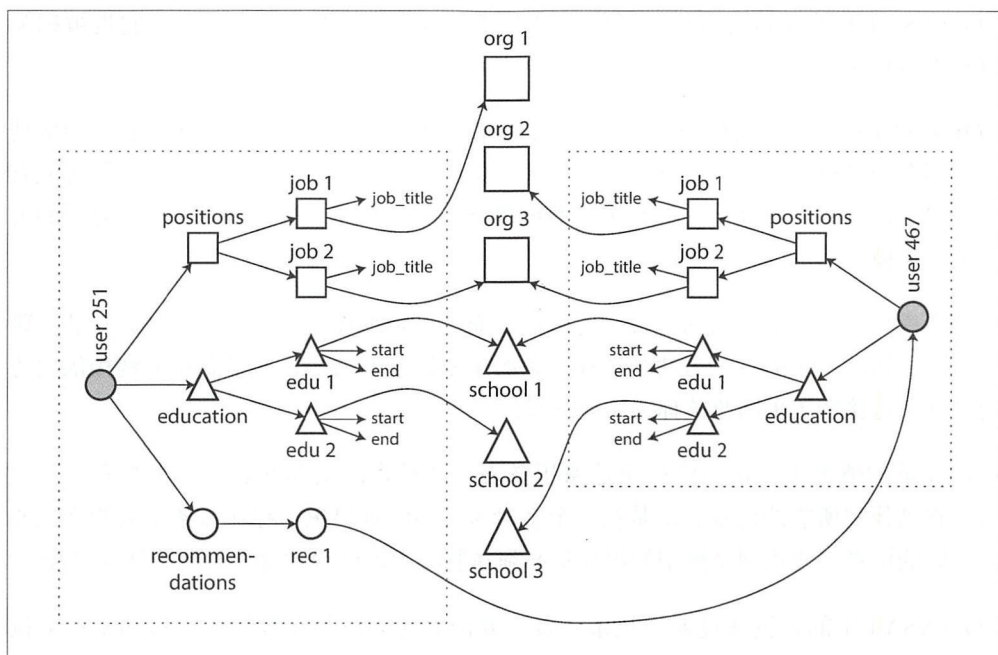


图2-4：采用多对多关系来扩展简历

IMS采用了相当简单的数据模型，称为层次模型，它与文档数据库使用的JSON模型有一些显著的相似之处<sup>[2]</sup>。它将所有数据表示为嵌套在记录中的记录（树），与图2-2的JSON结构非常相似。

和文档数据库类似，IMS可以很好地支持一对多关系中，但是它支持多对多关系则有些困难，而且不支持联结。开发人员必须决定是复制（反规范化）多份数据，还是手动解析记录之间的引用。20世纪六七十年代的这些问题与开发人员今天遇到的文档数据库的问题非常相似<sup>[15]</sup>。

为了解决层次模型的局限性，之后又提出了多种解决方案。其中最著名的是关系模型（relational model, 后来演变成为SQL，并被广泛接受）和网络模型（network model, 最初有很多拥趸，可惜最终被人们淡忘）。在20世纪70年代，这两个阵营之间的“大辩论”持续了很长时间<sup>[2]</sup>。

由于这两种模式所解决的问题与现在面临的争论非常相关，因此值得从现在的角度重新审视这场辩论。

## 网络模型

网络模型由一个称为数据系统语言会议（Conference on Data System Languages,



CODASYL)的委员会进行标准化,并由多个不同的数据库厂商实施,它也被称为CODASYL模型<sup>[16]</sup>。

CODASYL模型是层次模型的推广。在层次模型的树结构中,每个记录只有一个父结点;而在网络模型中,一个记录可能有多个父结点。例如,“大西雅图地区”可能有一个记录,居住在该地区的每个用户都链接指向它。从而支持对多对一和多对多的关系进行建模。

在网络模型中,记录之间的链接不是外键,而更像是编程语言中的指针(会存储在磁盘上)。访问记录的唯一方法是选择一条始于根记录的路径,并沿着相关链接依次访问。这条链接链条也因此被称为访问路径。

在最简单的情况下,访问路径像是遍历链表:从链表的头部开始,一次查看一个记录,直到找到所需的记录。但是在一个多对多关系的世界里,存在多条不同的路径通向相同的记录,使用网络模型的程序员必须在脑海中设法跟踪这些不同的访问路径。

CODASYL中的查询通过遍历记录列表,并沿着访问路径在数据库中移动光标来执行。如果记录有多个父结点(即来自其他记录的多个传入指针),则应用程序代码必须跟踪所有关系。甚至CODASYL委员会的成员也承认,这就像在一个 $n$ 维数据空间中进行遍历<sup>[17]</sup>。

在20世纪70年代,尽管手动路径选择能够最有效地利用当时非常有限的硬件资源(例如磁带驱动器,搜索速度非常慢),但最大的问题在于它们使查询和更新数据库变得异常复杂而没有灵活性。无论是层次模型还是网络模型,如果脱离数据的访问路径,那么将寸步难行。它也支持改变访问路径,但随之需要大量的手写数据库查询代码,重新实现处理新的访问路径。总之,对应用程序的数据模型进行更改是一件非常困难的事情。

## 关系模型

相比之下,关系模型所做的则是定义了所有数据的格式:关系(表)只是元组(行)的集合,仅此而已。没有复杂的嵌套结构,也没有复杂的访问路径。可以读取表中的任何一行或者所有行,支持任意条件查询。可以指定某些列作为键并匹配这些列来读取特定行。可以在任何表中插入新行,而不必担心与其他表之间的外键关系<sup>注4</sup>。

在关系数据库中,查询优化器自动决定以何种顺序执行查询,以及使用哪些索引。这

---

注4: 外键约束会限制某些修改,但关系模型不需要此类约束。即使有约束,在查询时也会执行外键联结,而在CODASYL中,联结则在插入时高效地完成。





些选择实际上等价于“访问路径”，但最大的区别在于它们是由查询优化器自动生成的，而不是由应用开发人员所维护，因此不用过多地考虑它们。

如果想用新的方式查询数据，只需声明一个新的索引，查询会自动使用最合适的索引。不需要更改查询即可利用新的索引（另请参阅本章后面的“数据查询语言”部分）。因此，关系模型使得应用程序添加新功能变得非常容易。

关系数据库的查询优化器称得上是一个复杂的怪兽，研究开发人员多年来持续投入，花费巨大<sup>[18]</sup>。不管怎样，关系模型的一个核心要点是：只需构建一次查询优化器，然后使用该数据库的所有应用程序都可以从中受益。如果没有查询优化器，那么为特定查询手动编写访问路径比编写通用优化器更容易，但从长远来看，通用解决方案更胜一筹。

## 文档数据库的比较

从以下角度来看，文档数据库是某种方式的层次模型：即在其父记录中保存了嵌套记录（一对多关系，如图2-1中的positions、education和contact\_info），而不是存储在单独的表中。

但是，在表示多对一和多对多的关系时，关系数据库和文档数据库并没有根本的不同：在这两种情况下，相关项都由唯一的标识符引用，该标识符在关系模型中被称为外键，在文档模型中被称为文档引用<sup>[9]</sup>。标识符可以查询时通过联结操作或相关后续查询来解析。迄今为止，文档数据库并未遵循CODASYL标准。

## 关系数据库与文档数据库现状

在比较关系数据库与文档数据库时，需要考虑很多方面的差异，包括它们的容错性（参阅第5章）和并发处理（参阅第7章）。本章将只关注数据模型中的差异。

支持文档数据模型的主要论点是模式灵活性，由于局部性而带来较好的性能，对于某些应用来说，它更接近于应用程序所使用的数据结构。关系模型则强在联结操作、多对一和多对多关系更简洁的表达上，与文档模型抗衡。

### 哪种数据模型的应用代码更简单？

如果应用数据具有类似文档的结构（即一对多关系树，通常一次加载整个树），那么使用文档模型更为合适。而关系型模型则倾向于某种数据分解，它把文档结构分解为多个表（如图2-1中的positions、education和contact\_info），有可能使得模式更为笨重，以及不必要的应用代码复杂化。



文档模型也有一定的局限性：例如，不能直接引用文档中的嵌套项，而需要说“用户251的职位列表中的第二项”（非常类似于层次模型中的访问路径）。然而，只要文档嵌套不太深，这通常不是问题。

在文档数据库中，对联结的支持不足是否是问题取决于应用程序。例如，在使用文档数据库记录事件发生时间的应用分析程序中，可能永远不需要多对多关系<sup>[19]</sup>。

但是，如果应用程序确实使用了多对多关系，那么文档模型就变得不太吸引人。可以通过反规范化来减少对联结的需求，但是应用程序代码需要做额外的工作来保持非规范化数据的一致性。通过向数据库发出多个请求，可以在应用程序代码中模拟联结，但是这也将会使应用程序变得复杂，并且通常比数据库内的专用代码执行的联结慢。在这些情况下，使用文档模型会导致应用程序代码更复杂、性能更差<sup>[15]</sup>。

通常无法一概而论哪种数据模型的应用代码更简单。这主要取决于数据项之间的关系类型。对于高度关联的数据，文档模型不太适合，关系模型可以胜任，而图模型（参阅本章后面的“图状数据模型”）则是最为自然的。

## 文档模型中的模式灵活性

大多数文档数据库，以及关系数据库中的JSON支持，都不会对文档中的数据强制执行任何模式。关系数据库中的XML通常支持带有可选的模式验证功能。没有模式意味着可以将任意的键-值添加到文档中，并且在读取时，客户端无法保证文档可能包含哪些字段。

文档数据库有时被称为无模式，但这具有误导性，因为读数据的代码通常采用某种结构因而存在某种隐形模式，而不是由数据库强制执行<sup>[20]</sup>。更准确的术语应该是读时模式（数据的结构是隐式的，只有在读取时才解释），与写时模式（关系数据库的一种传统方法，模式是显式的，并且数据库确保数据写入时都必须遵循）相对应<sup>[21]</sup>。

读时模式类似编程语言中的动态（运行时）类型检查，而写时模式类似于静态（编译时）类型检查。正如静态与动态类型检查的支持者对于它们的优缺点存在很大的争议一样<sup>[22]</sup>，数据库的模式执行也是一个有争议的话题，通常没有明确正确或错误的答案。

当应用程序需要改变数据格式时，这些方法之间的差异就变得尤其明显。例如，当前用户的全名存储在一个字段中，而现在想分别存储名字和姓氏<sup>[23]</sup>。在文档数据库中，只需使用新字段来编写新文档，并在应用层来处理读取旧文档的情况。例如：



```
if (user && user.name && !user.first_name) {  
    // 2013年12月8日之前写的文档，不存在first_name  
    user.first_name = user.name.split(" ")[0];  
}
```

而对于“静态类型”数据库模式中，通常会按照以下方式执行升级（migration）：

```
ALTER TABLE users ADD COLUMN first_name text;  
UPDATE users SET first_name = split_part(name, ' ', 1);      -- PostgreSQL  
UPDATE users SET first_name = substring_index(name, ' ', 1); -- MySQL
```

模式更改由于速度慢并且需要停机，因而评价不高。但这种坏名声其实并不太公平：大多数关系数据库系统可以在几毫秒内执行ALTER TABLE语句。MySQL则需要注意，它执行ALTER TABLE时会把现在的整张表复制，因而当表很大时可能会需要几分钟甚至几小时的停机时间，尽管现在有各种辅助工具可以解决这个限制<sup>[24-26]</sup>。

在大表上运行UPDATE语句，对于任何数据库都可能会很慢，因为每一行都需要重写。如果这是不可接受的，应用程序可以将first\_name设置为默认值NULL，并在读取时填充它，就像使用文档数据库一样。

如果集中的项由于某种原因（例如数据异构），并不都具有相同的结构，例如：

- 有许多不同类型的对象，将每种类型的对象都保存在各自的表中不太现实。
- 数据的结构由无法控制的外部系统所决定，而且可能随时改变。

在这些情况下，模式带来的损害大于它所能提供的帮助，无模式文档可能是更自然的数据模型。但是，当所有记录都有相同结构时，模式则是记录和确保这种结构的有效机制。我们将在第4章更详细地讨论模式和模式演变。

## 查询的数据局部性

文档通常存储为编码为JSON、XML或其二进制变体（如MongoDB的BSON）的连续字符串。如果应用程序需要频繁访问整个文档（例如，在网页上呈现它），则存储局部性具有性能优势。如果数据被划分在多个表中，如图2-1所示，则需要多次索引查找来检索所有数据，中间可能需要更多的磁盘I/O并花费更多的时间。

局部性优势仅适用需要同时访问文档大部分内容的场景。由于数据库通常会加载整个文档，如果应用只是访问其中的一小部分，则对于大型文档数据来讲就有些浪费。对文档进行更新时，通常会重写整个文档，而只有修改量不改变源文档大小时，原地覆盖更新才更有效<sup>[19]</sup>。因此，通常建议文档应该尽量小且避免写入时增加文档大小<sup>[9]</sup>。这些性能方面的不利因素大大限制了文档数据库的适用场景。





值得指出的是，将相关数据归为一组的局部性想法并不仅见于文档模型。例如，Google的Spanner数据库在关系数据模型中提供了相同的局部性，支持模式声明某些表的行应该在父表内交错（嵌套）<sup>[27]</sup>。Oracle支持类似的操作，称为“多表索引集群表”特性<sup>[28]</sup>，Bigtable数据模型（用于Cassandra和HBase）中的列族概念类似<sup>[29]</sup>。

我们将在第3章介绍更多局部性相关内容。

## 文档数据库与关系数据库的融合

自2000年中期以来，大多数关系数据库系统（MySQL除外）都支持XML。其中包括对XML文档进行本地修改，在XML文档中进行索引和查询等，这样应用程序可以获得与文档数据库非常相似的数据模型。

PostgreSQL版本9.3<sup>[8]</sup>、MySQL版本5.7以及IBM DB2 10.5<sup>[30]</sup>，都对JSON文档提供了相应支持。考虑到JSON在Web API中非常流行，其他关系数据库很可能将很快增加JSON支持。

文档数据库方面，RethinkDB的查询接口支持和关系型类似的联结，而一些MongoDB驱动程序可以自动解析数据库的引用关系（在客户端执行高效联结，注意，因为需要额外的网络往返，且优化程度较低，因此绝对性能可能慢于数据库端执行）。

随着时间的推移，似乎关系数据库与文档数据库变得越来越相近，或许这是一件好事：数据模型可以相互补充<sup>注5</sup>。如果数据库能够很好处理文档类数据，还能对其执行关系查询，那么应用程序可以使用最符合其需求的功能的组合。

融合关系模型与文档模型是未来数据库发展的一条很好的途径。

## 数据查询语言

当关系模型是初被引入时，就包含了查询数据的新方法：SQL是一种声明式查询语言，而IMS和CODASYL则是命令式。这种差别意味着什么呢？

很多常用的编程语言都是命令式。例如，如果有一个动物物种的列表，可能会写这样的代码来查询列表中的鲨鱼：

---

注5：Codd对关系模型的最早描述<sup>[1]</sup>实际上允许在关系模式中使用与JSON文档非常类似的东西，称之为非简单域。基本思想是，行中的值不一定只是数字或字符串一样的原始数据类型，也可以是嵌套关系（表），因此可以将任意嵌套的树结构作为值，很像30多年后SQL中支持JSON或XML。



```
function getSharks() {  
  var sharks = [];  
  for (var i = 0; i < animals.length; i++) {  
    if (animals[i].family === "Sharks") {  
      sharks.push(animals[i]);  
    }  
  }  
  return sharks;  
}
```

而对于关系代数，则会写成这样：

$$\text{sharks} = \sigma_{\text{family} = \text{"Sharks"}} (\text{animals})$$

其中  $\sigma$ （希腊字母  $\sigma$ ）是选择操作符，只返回符合条件的那些动物。

SQL遵循上述关系代数的结构：

```
SELECT * FROM animals WHERE family = 'Sharks';
```

命令式语言告诉计算机以特定顺序执行某些操作。你完全可以推理整个过程，逐行遍历代码、评估相关条件、更新对应的变量，并决定是否再循环一遍。

而对于声明式查询语言（如SQL或关系代数），则只需指定所需的数据模式，结果需满足什么条件，以及如何转换数据（例如，排序、分组和聚合），而不需指明如何实现这一目标。数据库系统的查询优化器会决定采用哪些索引和联结，以及用何种顺序来执行查询的各个语句。

声明式查询语言很有吸引力，它比命令式API更加简洁和容易使用。但更重要的是，它对外隐藏了数据库引擎的很多实现细节，这样数据库系统能够在不改变查询语句的情况下提高性能。

例如，在开头所示的命令式代码中，动物列表以特定顺序显示。如果数据库想要在后台回收未使用的磁盘空间，则可能需要移动记录，从而改变动物显示的顺序。数据库能否在不中断查询的情况下安全地执行此类操作？

SQL示例则不保证任何特定的顺序，所以顺序改变与否并不重要。但如果查询是命令式代码编写的，那么数据库永远无法确定代码是否依赖于排序。SQL在功能上有更多限制的事实，也给数据库提供了更多自动优化的空间。

最后，声明式语言通常适合于并行执行。现在CPU主要通过增加核，而不是通过比之前更高的时钟频率<sup>[31]</sup>来提升速度。而命令式代码由于指定了特定的执行顺序，很难在多核和多台机器上并行化。声明式语言则对于并行执行更为友好，它们仅指定了结果



所满足的模式，而不指定如何得到结果的具体算法。所以如果可以的话，数据库都倾向于采用并行方式实现查询语言<sup>[32]</sup>。

## Web上的声明式查询

声明式查询语言的优点不仅限于数据库。为了说明这一点，我们在另一个不同的环境中比较声明式和命令式的方法：Web浏览器。

假设有一个关于海洋动物的网站。用户当前正在查看有关鲨鱼的页面，因此用导航项“鲨鱼”标记为当前网页，如下所示：

```
<ul>
  <li class="selected"> ❶
    <p>Sharks</p> ❷
    <ul>
      <li>Great White Shark</li>
      <li>Tiger Shark</li>
      <li>Hammerhead Shark</li>
    </ul>
  </li>
  <li>
    <p>Whales</p>
    <ul>
      <li>Blue Whale</li>
      <li>Humpback Whale</li>
      <li>Fin Whale</li>
    </ul>
  </li>
</ul>
```

- ❶ 选定的项目被标记为CSS类别"selected"。
- ❷ <p>Sharks</p>表示当前页面的标题。

现在需要将当前页面的标题设置蓝色背景，这样视觉上更突出。采用CSS样式表很简单：

```
li.selected > p {
  background-color: blue;
}
```

这里，CSS选择器li.selected > p声明了要用蓝色样式：即直接父元素是<li>元素且CSS类为selected的所有<p>元素。示例中的元素<p> Sharks </p>匹配此模式，但是<p>Whales</p>不匹配，因为它的父元素<li>缺少class="selected"。

如果使用XSL而不是CSS，可以这样来实现：





```
<xsl:template match="li[@class='selected']/p">
  <fo:block background-color="blue">
    <xsl:apply-templates/>
  </fo:block>
</xsl:template>
```

这里，XPath表达式`li [@ class = 'selected']/p`相当于上例中的CSS选择器`li.selected > p`。CSS和XSL的相同之处在于，它们都是用于指定文档样式的声明式语言。

试想如果不得不采用命令式的方法，情况会是怎样？例如JavaScript使用核心文档对象模型（Document Object Model, DOM）API，最后结果可能是这样：

```
var liElements = document.getElementsByTagName("li");
for (var i = 0; i < liElements.length; i++) {
  if (liElements[i].className === "selected") {
    var children = liElements[i].childNodes;
    for (var j = 0; j < children.length; j++) {
      var child = children[j];
      if (child.nodeType === Node.ELEMENT_NODE && child.tagName === "P") {
        child.setAttribute("style", "background-color: blue");
      }
    }
  }
}
```

这段JavaScript代码命令式地将元素`<p>Sharks</p>` 设置为蓝色背景，但是代码很糟糕。它不仅比CSS和XSL等效版本更长、更难理解，而且还有一些严重的问题：

- 如果`selected`类被删除（例如，因为用户单击了不同的页面），那么即使代码重新运行，蓝色也不会被移除，因此该项将始终保持高亮显示，直到整个页面被重新加载。而使用CSS，浏览器会自动检测`li.selected > p`规则何时不再适用，并在`selected`类被删除后立即清除蓝色背景。
- 如果想利用新的API，比如`document.getElementsByClassName("selected")`，或者`document.evaluate()`，可能会提高性能，但是必须重写代码。另一方面，浏览器厂商可以在不破坏兼容性的情况下提高CSS和XPath的性能。

对于Web浏览器的例子，使用声明式CSS样式表比用JavaScript命令式地操作样式好得多。类似地，在数据库中，像SQL这样的声明式查询语言比命令式查询APIs要好得多<sup>注6</sup>。

---

注6： IMS和CODASYL都使用命令式查询API。应用程序通常使用COBOL代码遍历数据库中的记录，每次处理一条记录<sup>[2,16]</sup>。



## MapReduce查询

MapReduce是一种编程模型，用于在许多机器上批量处理海量数据，兴起于Google<sup>[33]</sup>。一些NoSQL存储系统（例如MongoDB和CouchDB）支持有限的MapReduce方式在大量文档上执行只读查询。

MapReduce在第10章中有更详细的描述。这里简要讨论一下MongoDB对该模型的使用。

MapReduce既不是声明式查询语言，也不是一个完全命令式的查询API，而是介于两者之间：查询的逻辑用代码片段来表示，这些代码片段可以被处理框架重复地调用。它主要基于许多函数式编程语言中的map（也称为collect）和reduce（也称为fold或inject）函数。

举个例子，假设你是一名海洋生物学家，每当你看到海洋中的动物时，就会在数据库中添加观察记录。现在你想生成一份报告，来说明每个月看到了多少鲨鱼。

在PostgreSQL中，可以这样表达该查询：

```
SELECT date_trunc('month', observation_timestamp) AS observation_month, ❶
       sum(num_animals) AS total_animals
FROM observations
WHERE family = 'Sharks'
GROUP BY observation_month;
```

- ❶ date\_trunc('month', timestamp) 函数用来获取时间戳中的月份，并返回该月份开始的另一个时间戳。换句话说，它将时间戳向下舍入到最近的月份。

此查询首先对观察结果进行过滤，仅显示鲨鱼家族的物种，然后按照他们发生的月份对观察结果进行分组，最后将该月的所有观察的动物数量求和汇总。

MongoDB中的MapReduce功能也可以实现类似目的，如下所示：

```
db.observations.mapReduce(
  function map() { ❷
    var year = this.observationTimestamp.getFullYear();
    var month = this.observationTimestamp.getMonth() + 1;
    emit(year + "-" + month, this.numAnimals); ❸
  },
  function reduce(key, values) { ❹
    return Array.sum(values); ❺
  },
  {
    query: { family: "Sharks" }, ❶
    out: "monthlySharkReport" ❻
  }
);
```



```
    }  
  );
```

- ❶ 过滤器声明式地指定鲨鱼种类（这是MongoDB对MapReduce的特有扩展）。
- ❷ 对于每个匹配查询的文档，都会调用一次JavaScript的map函数，并将其设置为文档对象。
- ❸ map函数发射一个键-值对，其中键是由年份和月份组成的字符串，如"2013-12"或"2014-1"；值代表观察的动物数量。
- ❹ map函数发射的键-值对按键分组。对于相同键（即相同的月份和年份）的所有键-值对，调用reduce函数。
- ❺ reduce函数将特定月份内所有观察到的动物数量相加。
- ❻ 最终的输出写入到monthlySharkReport集合中。

例如，假设observations集合包含如下两个文档：

```
{  
  observationTimestamp: Date.parse("Mon, 25 Dec 1995 12:34:56 GMT"),  
  family:      "Sharks",  
  species:     "Carcharodon carcharias",  
  numAnimals:  3  
}  
{  
  observationTimestamp: Date.parse("Tue, 12 Dec 1995 16:17:18 GMT"),  
  family:      "Sharks",  
  species:     "Carcharias taurus",  
  numAnimals:  4  
}
```

每个文档都会调用一次map函数，从而产生emit ("1995-12", 3) 和 ("1995-12", 4)。随后，reduce函数将被调用，reduce ("1995-12", [3,4])，返回7。

map和reduce函数对于可执行的操作有所限制。它们必须是纯函数，这意味着只能使用传递进去的数据作为输入，而不能执行额外的数据库查询，也不能有任何副作用。这些限制使得数据库能够在任何位置、以任意顺序来运行函数，并在失败时重新运行这些函数。不管怎样，该功能非常强大，可以通过它来解析字符串、调用库函数、执行计算等。

MapReduce是一个相当底层的编程模型，用于在计算集群上分布执行。而SQL这样的更高层次的查询语言可以通过一些MapReduce操作pipeline来实现（参阅第10章），





当然也有很多SQL的分布式实现并不借助MapReduce。请注意，SQL并没有任何限制规定它只能在单个机器上运行，而MapReduce也并非垄断了分布式查询。

能够在查询中使用JavaScript代码是一种高级查询特性，它也不限于MapReduce，某些SQL数据库也可以扩展使用JavaScript函数<sup>[34]</sup>。

MapReduce的一个可用性问题是，必须编写两个密切协调的JavaScript函数，这通常比编写单个查询更难。此外，声明式查询语言为查询优化器提供了更多提高查询性能的机会。由于这些原因，MongoDB 2.2增加了称为聚合管道的声明式查询语言的支持<sup>[9]</sup>。采用这种语言，鲨鱼计数查询可以实现如下：

```
db.observations.aggregate([
  { $match: { family: "Sharks" } },
  { $group: {
    _id: {
      year: { $year: "$observationTimestamp" },
      month: { $month: "$observationTimestamp" }
    },
    totalAnimals: { $sum: "$numAnimals" }
  } }
]);
```

聚合管道在表达能力上相当于SQL的子集，但是它使用了基于JSON的语法，而不是SQL的英语句式语法。或许这种差异仅仅是个品味问题。然而这个故事的寓意是，NoSQL系统可能会发现自己意外地重新发明了SQL，尽管是伪装的。

## 图状数据模型

我们之前看到，多对多关系是不同数据模型之间的重要区别特征。如果数据大多是一对多关系（树结构数据）或者记录之间没有关系，那么文档模型是最合适的。

但是，如果多对多的关系在数据中很常见呢？关系模型能够处理简单的多对多关系，但是随着数据之间的关联越来越复杂，将数据建模转化为图模型会更加自然。

图由两种对象组成：顶点（也称为结点或实体）和边（也称为关系或弧）。很多数据可以建模为图。典型的例子包括：

### 社交网络

顶点是人，边指示哪些人彼此认识。

### Web图

顶点是网页，边表示与其他页面的HTML链接。



## 公路或铁路网

顶点是交叉路口，边表示他们之间的公路或铁路线。

有很多著名的算法可以在这些图上运行。例如，汽车导航系统搜索道路网中任意两点之间的最短路径，PageRank可以计算Web图上网页的流行度，从而确定搜索排名。

在刚才给出的示例中，图的顶点表示相同类型的事物（分别是人、网页或交叉路口）。然而，图并不局限于这样的同构数据，图更为强大的用途在于，提供了单个数据存储区中保存完全不同类型对象的一致性方式。例如，Facebook维护了一个包含许多不同类型的顶点与边的大图：顶点包括人、地点、事件、签到和用户的评论；边表示哪些人是彼此的朋友，签到发生在哪些位置，谁评论了哪个帖子，谁参与了哪个事件等<sup>[35]</sup>。

本节我们将使用图2-5所示的示例。它可以来自于社交网络或某个族谱数据库。该例子显示了两个人，分别来自爱达荷州的Lucy和来自法国波恩的Alain。他们结婚了，目前住在伦敦。

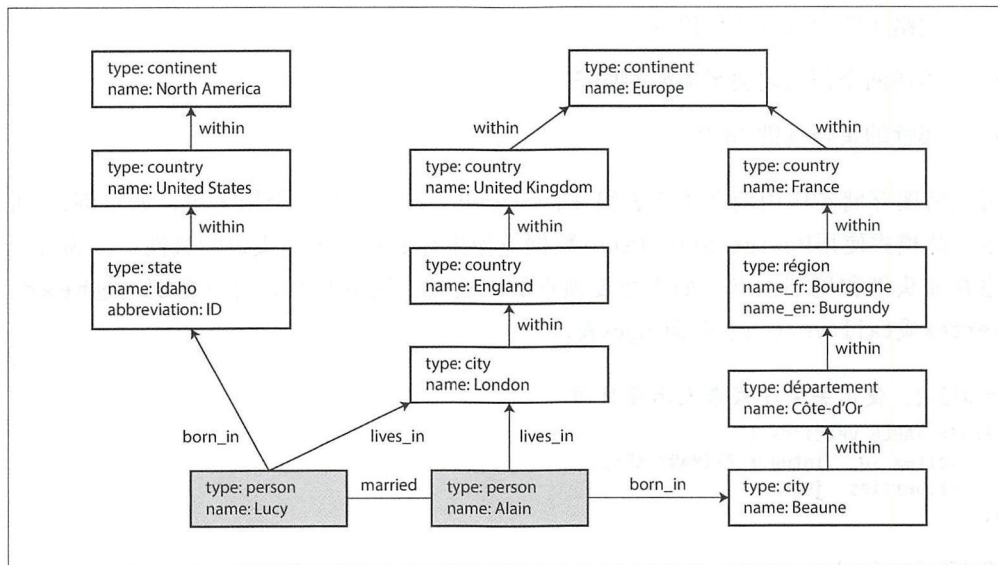


图2-5：图结构数据的例子（方框表示顶点，箭头表示边）

有多种不同但相关的方法可以构建和查询图中的数据。本节将讨论属性图模型（*property graph*，以Neo4j、Titan和InfiniteGraph为代表）和三元存储模型（*triple-store*，以Datomic、AllegroGraph等为代表）。我们将介绍三种声明式图查询语言：



Cypher、SPARQL和Datalog。之外，还有像Gremlin<sup>[36]</sup>这样的命令式图查询语言，以及Pregel这样的图处理框架（参阅第10章）。

## 属性图

在属性图模型中，每个顶点包括：

- 唯一的标识符。
- 出边的集合。
- 入边的集合。
- 属性的集合（键-值对）。

每个边包括：

- 唯一的标识符。
- 边开始的顶点（尾部顶点）。
- 边结束的顶点（头部顶点）。
- 描述两个顶点间关系类型的标签。
- 属性的集合（键-值对）。

可以将图存储看作由两个关系表组成，一个用于顶点，另一个用于边，如示例2-2所示（此模式使用PostgreSQL JSON数据类型来存储每个顶点或边的属性）。为每个边存储头部和尾部顶点，如果想要顶点的入边或出边的集合，可以分别通过`head_vertex`或`tail_vertex`来查询edges表。

示例2-2：使用关系模式来表示属性图

```
CREATE TABLE vertices (  
    vertex_id integer PRIMARY KEY,  
    properties json  
);  
  
CREATE TABLE edges (  
    edge_id integer PRIMARY KEY,  
    tail_vertex integer REFERENCES vertices (vertex_id),  
    head_vertex integer REFERENCES vertices (vertex_id),  
    label text,  
    properties json  
);  
  
CREATE INDEX edges_tails ON edges (tail_vertex);  
CREATE INDEX edges_heads ON edges (head_vertex);
```





关于图模型一些值得注意的地方：

1. 任何顶点都可以连接到其他任何顶点。没有模式限制哪种事物可以或不可以关联。
2. 给定某个顶点，可以高效地得到它的所有入边和出边，从而遍历图，即沿着这些顶点链条一直向前或向后（这就是为什么示例2-2中在tail\_vertex和head\_vertex列上都建立索引的原因）。
3. 通过对不同类型的关系使用不同的标签，可以在单个图中存储多种不同类型的信息，同时仍然保持整洁的数据模型。

这些特性为数据建模提供了很大的灵活性，如图2-5所示。图中显示了一些传统关系模式难以表达的东西，例如不同国家的不同类型的地区结构（法国有省和区，而美国有县和州），特殊历史原因，以及不同粒度的数据（Lucy当前住所被指定为一个城市，而她的出生地则是州一级）。

可以把这个图扩展到包括许多关于Lucy和Alain的其他信息，或者其他入。例如，可以用它来表示他们的任何食物过敏（通过为每个过敏源引入顶点，以及人与过敏原之间的边来表示过敏），并将过敏源与顶点的集合联结，这些顶点显示哪些食物含有哪些物质。然后，可以编写一个查询找出每个人吃什么是安全的。图有利于演化：向应用程序添加功能时，图可以容易地扩展以适应数据结构的不断变化。

## Cypher查询语言

Cypher是一种用于属性图的声明式查询语言，最早为Neo4j图形数据库而创建<sup>[37]</sup>（它以电影“黑客帝国”中的一个角色命名，与密码学中的密码无关<sup>[38]</sup>）。

示例2-3展示了将图2-5的左边部分插入图数据库的Cypher查询。图的其余部分可以类似添加，在此略过。每个顶点都有一个像USA或Idaho这样的符号名称，查询可以使用这些名称创建顶点之间的边，使用箭头符号：(Idaho) -[:WITHIN]->(USA)创建一个标签为WITHIN的边，其中Idaho为尾结点，USA为头结点。

示例2-3：对于图2-5中的一部分数据，采用Cypher方式查询

```
CREATE
  (NAmerica:Location {name:'North America', type:'continent'}),
  (USA:Location      {name:'United States', type:'country'  }),
  (Idaho:Location    {name:'Idaho',          type:'state'   }),
  (Lucy:Person       {name:'Lucy'  }),
  (Idaho) -[:WITHIN]-> (USA) -[:WITHIN]-> (NAmerica),
  (Lucy)  -[:BORN_IN]-> (Idaho)
```

把图2-5中所有的顶点和边插入到数据库之后，就可以解决一些很有意思的问题，例如



查找所有从美国移民到欧洲的人员名单。更准确地讲，我们需要查找BORN\_IN边指向美国，而LIVING\_IN边指向欧洲的所有的顶点，然后返回每个这样顶点的name属性。

示例2-4给出了如何用Cypher来实现这样的查询。MATCH语句中采用了相同的箭头语义(person) -[:BORN\_IN] -> ()来匹配这样的模式，即所有顶点间带有标签BORN\_IN的边，且尾部顶点对应于变量person，而头部顶点则没有要求。

示例2-4：采用Cypher查询从美国移民到欧洲的人员名单

```
MATCH
  (person) -[:BORN_IN]-> () -[:WITHIN*0..]-> (us:Location {name:'United States'}),
  (person) -[:LIVES_IN]-> () -[:WITHIN*0..]-> (eu:Location {name:'Europe'})
RETURN person.name
```

该查询的具体解读如下：

找到满足以下两个条件的任何顶点（顶点代表人，称其为person）：

1. person有一个到其他顶点的出边BORN\_IN。从该顶点开始，可以沿着一系列出边WITHIN，直到最终到达类型为Location的顶点，name属性为"United States"。
2. 同一个person顶点也有一个出边LIVES\_IN。沿着这条边，然后是一系列出边WITHIN，最终到达类型为Location的顶点，name属性为"Europe"。

对于每个这样的顶点，返回name属性。

其实有多种可行的方法执行上述查询。这里给出的方案是，首先扫描数据库中的所有入，检查每个人的出生地和居住地，然后只返回符合条件的人。

当然也可以采用其他等价方式，例如可以从两个Location顶点开始。如果name属性上有索引，则可以高效地找到代表美国和欧洲的两个顶点。然后，通过沿着所有的WITHIN入边，可以继续查找美国和欧洲的所有地点（州、地区、城市等）。最后，根据某个位置顶点的入边BORN\_IN或LIVES\_IN来找到符合条件的人员。

对于声明式查询语言，通常在编写查询语句时，不需要指定执行细节：查询优化器会自动选择效率最高的执行策略，因此开发者可以专注于应用的其他部分。

## SQL中的图查询

示例2-2表明可以采用关系数据库表示图数据。如果把图数据放在关系结构中，是否意味着也可以支持SQL查询呢？

答案是肯定的，但存在一些困难。在关系数据库中，通常会预先知道查询中需要哪些



join操作。而对于图查询，在找到要查找的顶点之前，可能需要遍历数量未知的边，也就是说，join操作数量并不是预先确定的。

对于示例2-4，这主要集中在Cypher查询的()`-[:WITHIN * 0 ..]` `->` ()规则中。一个人的LIVES\_IN边可以指向任何类型的位置，例如街道、城市、地区、区域、国家等。城市可以位于(WIHING)某地区内，地区可以位于(WIHING)某州内，州位于(WIHING)国家内等。LIVES\_IN边可以直接指向正在查找的位置顶点，也可以是位置层次结构中删除的某些层。

Cypher可以用:WITHIN\*0.. 非常简洁地表达这个情况：它表示“沿着一个WITHIN边，遍历零次或多次”，就像正则表达式中的\*运算符（表示匹配零次或多次）那样。

SQL:1999标准以后，查询过程中这种可变的遍历路径可以使用称为递归公用表表达式（即WITH RECURSIVE语法）来表示。示例2-5采用该技术的SQL表达来执行相同的查询（查找从美国移民到欧洲的人员名单），目前PostgreSQL、IBM DB2、Oracle和SQL Server等都支持该技术，但与Cypher相比，语法仍显得非常笨拙。

示例2-5：在SQL中采用递归公用表表达式来执行与示例2-4相同的查询  
WITH RECURSIVE

```
-- in_usa is the set of vertex IDs of all locations within the United States
in_usa(vertex_id) AS (
    SELECT vertex_id FROM vertices WHERE properties->>'name' = 'United States' ❶
    UNION
    SELECT edges.tail_vertex FROM edges ❷
    JOIN in_usa ON edges.head_vertex = in_usa.vertex_id
    WHERE edges.label = 'within'
),

-- in_europe is the set of vertex IDs of all locations within Europe
in_europe(vertex_id) AS (
    SELECT vertex_id FROM vertices WHERE properties->>'name' = 'Europe' ❸
    UNION
    SELECT edges.tail_vertex FROM edges
    JOIN in_europe ON edges.head_vertex = in_europe.vertex_id
    WHERE edges.label = 'within'
),

-- born_in_usa is the set of vertex IDs of all people born in the US
born_in_usa(vertex_id) AS ( ❹
    SELECT edges.tail_vertex FROM edges
    JOIN in_usa ON edges.head_vertex = in_usa.vertex_id
    WHERE edges.label = 'born_in'
),

-- lives_in_europe is the set of vertex IDs of all people living in Europe
lives_in_europe(vertex_id) AS ( ❺
```





```
SELECT edges.tail_vertex FROM edges
  JOIN in_europe ON edges.head_vertex = in_europe.vertex_id
 WHERE edges.label = 'lives_in'
)

SELECT vertices.properties->>'name'
FROM vertices
-- join to find those people who were both born in the US *and* live in Europe
JOIN born_in_usa      ON vertices.vertex_id = born_in_usa.vertex_id ❹
JOIN lives_in_europe ON vertices.vertex_id = lives_in_europe.vertex_id;
```

- ❶ 首先找到name属性值为"United States"的顶点，并将其作为顶点集in\_usa中的第一个元素。
- ❷ 沿着集合in\_usa中顶点的所有入边within，并将它们添加到同一集合，直到遍历所有的入边。
- ❸ 从name属性值为"Europe"的顶点开始执行同样的操作，并建立顶点集in\_europe。
- ❹ 对于in\_usa集合中的每个顶点，按照入边born\_in来查找出生在美国境内某个地方的人。
- ❺ 类似地，对于in\_europe集合中的每个顶点，按照入边lives\_in来查找居住在欧洲的人。
- ❻ 最后，通过join把在美国出生的人的集合与在欧洲居住的人的集合相交。

如果相同的查询可以用一种查询语言写4行代码完成，而另一种查询语言则需要29行代码，这足以说明不同的数据模型适用于不同的场景。因此，选择适合应用程序的数据模型非常重要。

## 三元存储与SPARQL

三元存储模式几乎等同于属性图模型，只是使用不同的名词描述了相同的思想。尽管如此，考虑到有多种针对三元存储的工具和语言，它们可能是构建应用程序宝贵的补充，因此还是值得在此讨论。

在三元存储中，所有信息都以非常简单的三部分形式存储（主体，谓语，客体）。例如，在三元组（吉姆，喜欢，香蕉）中，吉姆是主体，喜欢是谓语（动词），香蕉是客体。

三元组的主体相当于图中的顶点。而客体则是以下两种之一：

1. 原始数据类型中的值，如字符串或数字。在这种情况下，三元组的谓语和客体分



别相当于主体（顶点）属性中的键和值。例如， $(lucy, age, 33)$ 就好比是顶点lucy，具有属性{"age": 33}。

2. 图中的另一个顶点。此时，谓语是图中的边，主体是尾部顶点，而客体是头部顶点。例如，在 $(lucy, marriedTo, alain)$ 中，主体lucy和客体alain都是顶点，并且谓语marriedTo是连接二者的边的标签。

示例2-6表示了与示例2-3中相同的数据，它以三元组的格式定义了Turtle，Turtle可以认为是Notation3(N3)的一个子集<sup>[39]</sup>。

示例2-6：以Turtle三元组方式表示图2-5中一部分数据

```
@prefix : <urn:example:>.
_:lucy   a      :Person.
_:lucy   :name   "Lucy".
_:lucy   :bornIn _:idaho.
_:idaho  a      :Location.
_:idaho  :name   "Idaho".
_:idaho  :type   "state".
_:idaho  :within _:usa.
_:usa    a      :Location.
_:usa    :name   "United States".
_:usa    :type   "country".
_:usa    :within _:namerica.
_:namerica a    :Location.
_:namerica :name  "North America".
_:namerica :type  "continent".
```

在上述示例中，图的顶点被写为\_*someName*。顶点的名字在定义文件以外没有任何意义，只是为了区分三元组的不同顶点。谓语表示边时，客体是另一个顶点，如\_*idaho* :within *\_:usa*。当谓语表示一个属性时，该客体则是一个字符串，如\_*usa* :name "United States"。

如果要定义相同主体的多个三元组，反复输入相同的单词就略显枯燥。可以使用分号来说明同一主体的多个对象信息。这样Turtle格式就更为简洁、可读性强，参见示例2-7。

示例2-7：以更简洁的语法重写示例2-6

```
@prefix : <urn:example:>.
_:lucy   a :Person; :name "Lucy";           :bornIn _:idaho.
_:idaho  a :Location; :name "Idaho";         :type "state"; :within _:usa.
_:usa    a :Location; :name "United States"; :type "country"; :within _:namerica.
_:namerica a :Location; :name "North America"; :type "continent".
```

## 语义网

如果阅读更多关于三元存储的信息，很可能会被卷入关于语义网大量文章漩涡之中。



三元存储数据模型其实完全独立于语义网，例如，Datomic<sup>[40]</sup>是一个三元存储，它与语义网并没有任何关系<sup>注7</sup>。然而，考虑到很多人认为两者紧密相连，有必要在这里做一些简单澄清。

语义网从本质上讲源于一个简单而合理的想法：网站通常将信息以文字和图片方式发布给人类阅读，那为什么不把信息发布为机器可读的格式给计算机阅读呢？资源描述框架（Resource Description Framework，RDF<sup>[41]</sup>就是这样一种机制，它让不同网站以一致的格式发布数据，这样来自不同网站的数据自动合并成一个数据网络，一种互联网级别包含所有数据的数据库。

不幸的是，语义网在21世纪初被严重夸大了，时至今日也没有在实践中见到任何靠谱的实现，由此许多人开始怀疑它。另外，还有其他一些方面的批评，包括令人眼花缭乱的各种缩略词、极其复杂的标准提议，以及过于自大的标榜。

## RDF数据模型

示例2-7中使用的Turtle语言代表了RDF数据的人类可读格式。有时候，RDF也用XML格式编写，不过更冗长一些，参看示例2-8。人眼更容易阅读像Turtle/N3这种格式，而采用Apache Jena<sup>[42]</sup>等自动化工具可以快速转换不同的RDF格式。

示例2-8：用RDF/XML语法表示示例2-7的数据

```
<rdf:RDF xmlns="urn:example:"  
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#">  
  
  <Location rdf:nodeID="idaho">  
    <name>Idaho</name>  
    <type>state</type>  
    <within>  
      <Location rdf:nodeID="usa">  
        <name>United States</name>  
        <type>country</type>  
        <within>  
          <Location rdf:nodeID="namerica">  
            <name>North America</name>  
            <type>continent</type>  
          </Location>  
        </within>  
      </Location>  
    </within>  
  </Location>  
  
  <Person rdf:nodeID="lucy">  
    <name>Lucy</name>  
    <bornIn rdf:nodeID="idaho"/>  
  </Person>  
</rdf:RDF>
```

注7：从技术上讲，Datomic使用五元组而不是三元组；另外两个字段是用于标识多版本的元数据。





```
</Person>
</rdf:RDF>
```

因为旨在为全网数据交换而设计，RDF存在一些特殊的约定。三元组的主体、谓语和客体通常是URI。例如谓语可能是URI，诸如<http://my-company.com/namespace#within>或<http://my-company.com/namespace#lives\_in>，而不仅仅是WITHIN或LIVES\_IN。这种设计背后的原因是，它假设你的数据需要和其他人的数据相结合，万一不同人给单词within或者lives\_in附加了不同的含义，采用URI则可以避免冲突，即谓语实际上是<http://other.org/foo#within>和<http://other.org/foo#lives\_in>。

从RDF的角度来看，URL <http://my-company.com/namespace>不一定需要解析出特定的内容，它更多的只是一个命名空间。为避免与http://URL混淆，上述示例使用了不可解析的URI，如urn:example:within。还好，只需在文件头部指定一次前缀，就可以忽略它。

## SPARQL查询语言

SPARQL是一种采用RDF数据模型的三元存储查询语言<sup>[43]</sup>，名字是SPARQL Protocol和RDF Query Language的缩写，发音为“sparkle”。它比Cypher更早，并且由于Cypher的模式匹配是借用SPARQL的，所以二者看起来非常相似<sup>[37]</sup>。

执行同样的查询（从美国移民到欧洲的人员），SPARQL比Cypher更加简洁（参看示例2-9）。

示例2-9：采用SPARQL来实现与示例2-4相同的查询

PREFIX : <urn:example:>

```
SELECT ?personName WHERE {
  ?person :name ?personName.
  ?person :bornIn / :within* / :name "United States".
  ?person :livesIn / :within* / :name "Europe".
}
```

可以看到总体结构非常相似。以下两个表达式是等价的（SPARQL中变量以问号开头）：

```
(person) -[:BORN_IN]-> () -[:WITHIN*0..]-> (location) # Cypher
?person :bornIn / :within* ?location. # SPARQL
```

由于RDF不区分属性和边，可以同时对两者执行谓语操作，所以可以采用相同的语法来匹配属性上的查询条件。在下面的表达式中，变量usa被绑定到所有具有name属



性，且name属性值为字符串"United States"的顶点：

```
(usa {name:'United States'}) # Cypher
```

```
?usa :name "United States". # SPARQL
```

SPARQL是一种非常优秀的查询语言，即使语义网从未实际出现，它也可以成为应用程序内部使用的强大查询工具。

## 图数据库与网络模型的比较

在前面的“文档数据库是否在重演历史”中，我们介绍了CODASYL与关系模型如何竞争解决IMS中的多对多关系问题。乍一看，CODASYL的网络模型和图模型很相似，那么图形数据库是乔装打扮的另一个CODASYL吗？

答案是否定的，它们在以下几个重要方面有着明显区别：

- 在CODASYL中，数据库有一个模式来指定哪种记录类型可以嵌套在其他记录类型中。在图数据库中，则没有这样的限制：任何顶点都可以有边连接其他任何顶点。这就为应用程序适应不断变化的需求提供了更大的灵活性。
- 在CODASYL中，获取特定记录的唯一方法是遍历其中的一条访问路径。在图数据库中，则可以通过顶点的唯一ID直接引用该顶点，也可以使用索引查找满足特定值的那些顶点。
- 在CODASYL中，记录的子记录是有序集合，所以数据库必须保持这种排序（这对存储布局产生影响），当应用插入新记录时不得不考虑新记录在这些集合中的位置。在图数据库中，顶点和边不是有序的（只能在查询时对结果进行排序）。
- 在CODASYL中，所有的查询都是命令式的，难以编写，并且很容易被模式的变化所破坏。在图数据库中，可以采用命令式代码来实现自己的遍历，但大多数图形数据库还支持高级声明式查询语言，例如Cypher或SPARQL。

## Datalog基础

Datalog是比SPARQL或Cypher更为古老的语言，在20世纪80年代被学者广泛研究<sup>[44-46]</sup>。虽然在软件工程师中知名度较低，但它为以后的查询语言奠定了基础，因此它非常重要。



实践中有几个数据库系统采用了Datalog。例如它是Datomic<sup>[40]</sup>系统的查询语言，而Cascalog<sup>[47]</sup>是用于查询Hadoop大数据集的Datalog实现<sup>注8</sup>。

Datalog的数据模型类似于三元存储模式，但更为通用一些。它采用“谓语（主体，客体）”的表达方式而不是三元组（主体，谓语，客体）。示例2-10展示了如何使用Datalog。

示例2-10：采用Datalog表示图2-5中的数据子集

```
name(america, 'North America').
type(america, continent).
```

```
name(usa, 'United States').
type(usa, country).
within(usa, america).
```

```
name(idaho, 'Idaho').
type(idaho, state).
within(idaho, usa).
```

```
name(lucy, 'Lucy').
born_in(lucy, idaho).
```

已经定义好了数据之后，可以执行之前类似的查询，如示例2-11所示。它看起来有点不同于Cypher或SPARQL，但不要让这种情况使你放弃。Datalog是Prolog的子集，如果你学过相关计算机课程，或许你之前曾见过Prolong。

示例2-11：采用Datalog实现示例2-4相同的查询

```
within_recursive(Location, Name) :- name(Location, Name).      /* 规则 1 */
```

```
within_recursive(Location, Name) :- within(Location, Via),      /* 规则 2 */
                                     within_recursive(Via, Name).
```

```
migrated(Name, BornIn, LivingIn) :- name(Person, Name),        /* 规则 3 */
                                     born_in(Person, BornLoc),
                                     within_recursive(BornLoc, BornIn),
                                     lives_in(Person, LivingLoc),
                                     within_recursive(LivingLoc, LivingIn).
```

```
?- migrated(Who, 'United States', 'Europe').
/* Who = 'Lucy'. */
```

Cypher和SPARQL通过类似SELECT一次完成查询，而Datalog则每次实现一块。我们定义了告诉数据库关于新谓语的规则，例如两个新的谓语within\_recursive和migrated。这些谓语并不是存储在数据库中的三元组，而是从数据或其他规则派生而

---

注8： Datomic和Cascalog使用Datalog的Clojure S表达式语法。在下面的例子中，我们使用了更容易阅读的Prolog语法，但是功能方面没有差别。





来。规则可以引用其他规则，就像函数可以调用其他函数或者递归调用自己一样。像这样，复杂的查询可以通过每次完成一小块而逐步构建。

在规则中，以大写字母开头的单词是变量，谓词匹配与Cypher和SPARQL一样。例如，`name(Location, Name)`与使用了变量绑定`Location=namerica`和`Name='North America'`的三元组`name(namerica, 'North America')`等价。

如果系统可以在操作符`-`的右侧找到与所有谓词的匹配项，规则适用。当规则适用时，就好像将`-`的左侧添加到数据库中一样（其中的变量替换为它们匹配的值）。

因此，应用该规则的一种可能的工作步骤如下所示：

1. 数据库中存在`name(namerica, 'North America')`，因此规则1适用。它产生`within_recursive(namerica, 'North America')`。
2. 数据库中存在`within(usa, namerica)`，并且步骤1生成了`within_recursive(namerica, 'North America')`，所以规则2适用。它会产生`within_recursive(usa, 'North America')`。
3. 数据库中存在`within(idaho, usa)`，并且步骤2生成了`within_recursive(usa, 'North America')`，所以规则2适用。它产生`within_recursive(idaho, 'North America')`。

通过反复应用规则1和规则2，`within_recursive`谓词可以返回数据库中包含的所有North America地点（或任何其他地点名称）。该过程如图2-6所示。

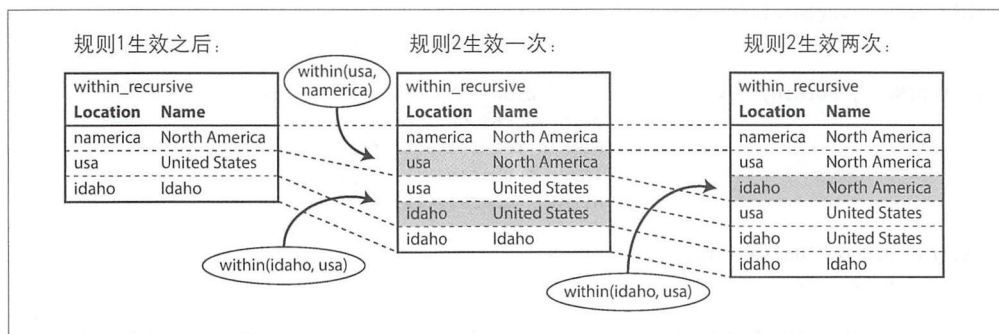


图2-6：使用示例2-11中的Datalog规则确定爱达荷州在北美

现在规则3可以找到出生在某个地方（BornIn）且居住在某个地方（LivingIn）的所有人。例如通过查询`BornIn='United States'`和`LivingIn='Europe'`，并将此人作为变量Who，要求Datalog系统找出变量Who包含哪些值。最终我们得到了与前面Cypher



和SPARQL查询相同的结果。

Datalog方法需要采取与其他查询语言略有不同的思维方式，但它非常强大，特别是规则可以在不同的查询中组合和重用。对于简单的一次性查询来说，这或许不太方便，但是如果数据非常复杂，处理起来会更加游刃有余。

## 小结

数据模型是一个庞大的主题，本章我们快速介绍了各种不同的数据模型。因篇幅所限，我们无法详细介绍每个模型的细节，但是希望这个概述足以激起你的兴趣，进一步来寻找最适合应用需求的模型。

历史上，数据最初被表示为一棵大树（层次模型），但是这不利于表示多对多关系，所以发明了关系模型来解决这个问题。最近，开发人员发现一些应用程序也不太适合关系模型。新的非关系“NoSQL”数据存储的两个主要方向上存在分歧：

1. 文档数据库的目标用例是数据来自于自包含文档，且一个文档与其他文档之间的关联很少。
2. 图数据库则针对相反的场景，目标用例是所有数据都可能会互相关联。

所有这三种模型（文档模型、关系模型和图模型），如今都有广泛使用，并且在各自的目标领域都足够优秀。我们观察到，一个模型可以用另一个模型来模拟。例如，图数据可以在关系数据库中表示，虽然处理起来比较笨拙。这就是为什么不同的系统用于不同的目的，而不是一个万能的解决方案。

文档数据库和图数据库有一个共同点，那就是它们通常不会对存储的数据强加某个模式，这可以使应用程序更容易适应不断变化的需求。但是，应用程序很可能仍然假定数据具有一定的结构，只不过是模式是显式（写时强制）还是隐式（读时处理）的问题。

每个数据模型都有自己的查询语言或框架，我们讨论了几个例子：SQL、MapReduce、MongoDB的聚合管道、Cypher、SPARQL和Datalog。我们还讨论了CSS和XSL/XPath，它们并不属于数据库查询语言，但存在有趣的相似之处。

虽然已经覆盖了很广的范围，但仍然有一些数据模型尚未提及。举几个简单的例子：

- 使用基因组数据的研究人员经常需要执行序列相似性搜索，这意味着需要用一个非常长的字符串（代表一个DNA分子），与存在相似但却不完全相同的大型字



符串数据库进行匹配。以上介绍的所有数据库都不适用于这种场景，这就是为什么研究人员开发了专门的基因组数据库软件，如GenBank<sup>[48]</sup>。

- 数十年来，粒子物理学家一直在进行海量数据的超大规模数据分析，像大型强子对撞机（LHC）这样的项目，现在可以处理数百PB级别的数据！在这种规模下，需要一些定制解决方案来避免硬件成本失控<sup>[49]</sup>。
- 全文搜索可以说是一种经常与数据库一起使用的数据库模型。信息检索是一个很大的专业课题，本书不会详细介绍，但是将在第3章和第三部分中介绍搜索索引相关内容。

本章暂时告一段路。在下一章中，我们将讨论在实现本章所描述的数据模型过程中有哪些重要的权衡设计。

## 参考文献

[1] Edgar F. Codd: “A Relational Model of Data for Large Shared Data Banks,” *Communications of the ACM*, volume 13, number 6, pages 377-387, June 1970. doi: 10.1145/362384.362685.

[2] Michael Stonebraker and Joseph M. Hellerstein: “What Goes Around Comes Around,” in *Readings in Database Systems*, 4th edition, MIT Press, pages 2-41, 2005. ISBN: 978-0-262-69314-1.

[3] Pramod J. Sadalage and Martin Fowler: *NoSQL Distilled*. Addison-Wesley, August 2012. ISBN: 978-0-321-82662-6.

[4] Eric Evans: “NoSQL: What’s in a Name?,” *blog.sym-link.com*, October 30, 2009.

[5] James Phillips: “Surprises in Our NoSQL Adoption Survey,” *blog.couchbase.com*, February 8, 2012.

[6] Michael Wagner: *SQL/XML:2006 - Evaluierung der Standardkonformität ausgewählter Datenbanksysteme*. Diplomica Verlag, Hamburg, 2010. ISBN: 978-3-836-64609-3.

[7] “XML Data in SQL Server,” SQL Server 2012 documentation, *technet.microsoft.com*, 2013.

[8] “PostgreSQL 9.3.1 Documentation,” The PostgreSQL Global Development Group, 2013.





- [9] “The MongoDB 2.4 Manual,” MongoDB, Inc., 2013.
- [10] “RethinkDB 1.11 Documentation,” *rethinkdb.com*, 2013.
- [11] “Apache CouchDB 1.6 Documentation,” *docs.couchdb.org*, 2014.
- [12] Lin Qiao, Kapil Surlaker, Shirshanka Das, et al.: “On Brewing Fresh Espresso: LinkedIn’s Distributed Data Serving Platform,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2013.
- [13] Rick Long, Mark Harrington, Robert Hain, and Geoff Nicholls: *IMS Primer*. IBM Redbook SG24-5352-00, IBM International Technical Support Organization, January 2000.
- [14] Stephen D. Bartlett: “IBM’s IMS—Myths, Realities, and Opportunities,” The Clipper Group Navigator, TCG2013015LI, July 2013.
- [15] Sarah Mei: “Why You Should Never Use MongoDB,” *sarahmei.com*, November 11, 2013.
- [16] J. S. Knowles and D. M. R. Bell: “The CODASYL Model,” in *Databases—Role and Structure: An Advanced Course*, edited by P. M. Stocker, P. M. D. Gray, and M. P. Atkinson, pages 19-56, Cambridge University Press, 1984. ISBN: 978-0-521-25430-4.
- [17] Charles W. Bachman: “The Programmer as Navigator,” *Communications of the ACM*, volume 16, number 11, pages 653-658, November 1973. doi: 10.1145/355611.362534.
- [18] Joseph M. Hellerstein, Michael Stonebraker, and James Hamilton: “Architecture of a Database System,” *Foundations and Trends in Databases*, volume 1, number 2, pages 141-259, November 2007. doi:10.1561/1900000002.
- [19] Sandeep Parikh and Kelly Stirman: “Schema Design for Time Series Data in MongoDB,” *blog.mongodb.org*, October 30, 2013.
- [20] Martin Fowler: “Schemaless Data Structures,” *martinfowler.com*, January 7, 2013.
- [21] Amr Awadallah: “Schema-on-Read vs. Schema-on-Write,” at *Berkeley EECS RAD Lab Retreat*, Santa Cruz, CA, May 2009.
- [22] Martin Odersky: “The Trouble with Types,” at *Strange Loop*, September 2013.



- [23] Conrad Irwin: “MongoDB—Confessions of a PostgreSQL Lover,” at *HTML5DevConf*, October 2013.
- [24] “Percona Toolkit Documentation: pt-online-schema-change,” Percona Ireland Ltd., 2013.
- [25] Rany Keddo, Tobias Bielohlawek, and Tobias Schmidt: “Large Hadron Migrator,” SoundCloud, 2013.
- [26] Shlomi Noach: “gh-ost: GitHub’s Online Schema Migration Tool for MySQL,” *githubengineering.com*, August 1, 2016.
- [27] James C. Corbett, Jeffrey Dean, Michael Epstein, et al.: “Spanner: Google’s Globally-Distributed Database,” at *10th USENIX Symposium on Operating System Design and Implementation (OSDI)*, October 2012.
- [28] Donald K. Burleson: “Reduce I/O with Oracle Cluster Tables,” *dba-oracle.com*.
- [29] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, et al.: “Bigtable: A Distributed Storage System for Structured Data,” at *7th USENIX Symposium on Operating System Design and Implementation (OSDI)*, November 2006.
- [30] Bobbie J. Cochrane and Kathy A. McKnight: “DB2 JSON Capabilities, Part 1: Introduction to DB2 JSON,” IBM developerWorks, June 20, 2013.
- [31] Herb Sutter: “The Free Lunch Is Over: A Fundamental Turn Toward Concurrency in Software,” *Dr. Dobbs’s Journal*, volume 30, number 3, pages 202-210, March 2005.
- [32] Joseph M. Hellerstein: “The Declarative Imperative: Experiences and Conjectures in Distributed Logic,” Electrical Engineering and Computer Sciences, University of California at Berkeley, Tech report UCB/EECS-2010-90, June 2010.
- [33] Jeffrey Dean and Sanjay Ghemawat: “MapReduce: Simplified Data Processing on Large Clusters,” at *6th USENIX Symposium on Operating System Design and Implementation (OSDI)*, December 2004.
- [34] Craig Kerstiens: “JavaScript in Your Postgres,” *blog.heroku.com*, June 5, 2013.
- [35] Nathan Bronson, Zach Amsden, George Cabrera, et al.: “TAO: Facebook’s Distributed Data Store for the Social Graph,” at *USENIX Annual Technical Conference (USENIX ATC)*, June 2013.

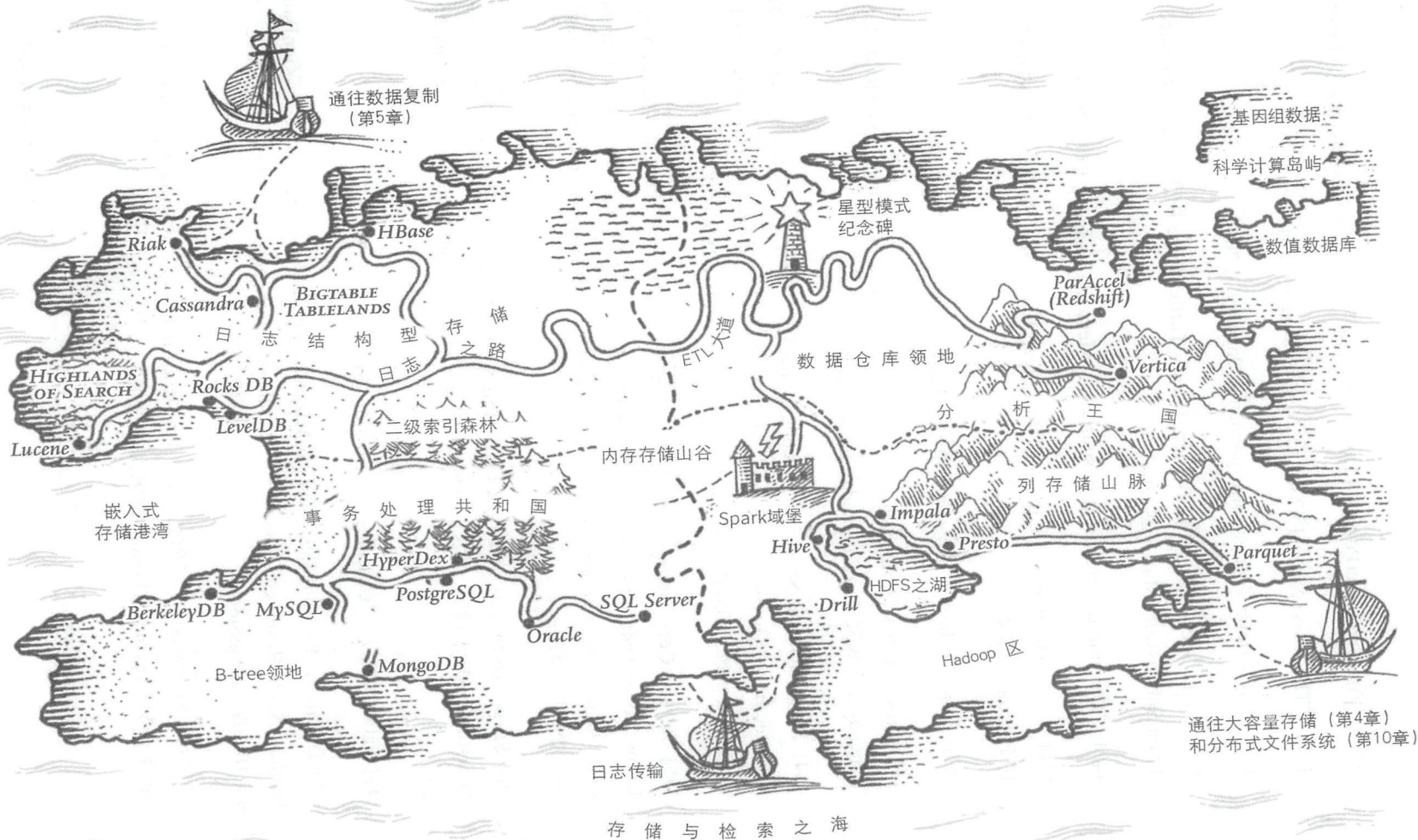


- [36] “Apache TinkerPop3.2.3 Documentation,” *tinkerpop.apache.org*, October 2016.
- [37] “The Neo4j Manual v2.0.0,” Neo Technology, 2013.
- [38] Emil Eifrem: Twitter correspondence, January 3, 2014.
- [39] David Beckett and Tim Berners-Lee: “Turtle–Terse RDF Triple Language,” W3C Team Submission, March 28, 2011.
- [40] “Datomic Development Resources,” Metadata Partners, LLC, 2013.
- [41] W3C RDF Working Group: “Resource Description Framework (RDF),” *w3.org*, 10 February 2004.
- [42] “Apache Jena,” Apache Software Foundation.
- [43] Steve Harris, Andy Seaborne, and Eric Prud’hommeaux: “SPARQL 1.1 Query Language,” W3C Recommendation, March 2013.
- [44] Todd J. Green, Shan Shan Huang, Boon Thau Loo, and Wenchao Zhou: “Datalog and Recursive Query Processing,” *Foundations and Trends in Databases*, volume 5, number 2, pages 105-195, November 2013. doi:10.1561/19000000017
- [45] Stefano Ceri, Georg Gottlob, and Letizia Tanca: “What You Always Wanted to Know About Datalog (And Never Dared to Ask),” *IEEE Transactions on Knowledge and Data Engineering*, volume 1, number 1, pages 146-166, March 1989. doi:10.1109/69.43410.
- [46] Serge Abiteboul, Richard Hull, and Victor Vianu: *Foundations of Databases*. Addison-Wesley, 1995. ISBN: 978-0-201-53771-0, available online at *webdam.inria.fr/Alice*
- [47] Nathan Marz: “Cascalog,” *cascalog.org*
- [48] Dennis A. Benson, Ilene Karsch-Mizrachi, David J. Lipman, et al.: “GenBank,” *Nucleic Acids Research*, volume 36, Database issue, pages D25-D30, December 2007. doi:10.1093/nar/gkm929
- [49] Fons Rademakers: “ROOT for Big Data Analysis,” at Workshop on the Future of Big Data Management, London, UK, June 2013.





# 分布式数据大洋





## 第3章

# 数据存储与检索

如果你把东西整理得井井有条，下次就不用查找了。

——德国谚语

从最基本的层面看，数据库只需做两件事情：向它插入数据时，它就保存数据；之后查询时，它应该返回那些数据。

在第2章中，我们讨论了数据模型和查询语言，即关于应用开发人员向数据库指明数据格式并在之后如何查询的机制。本章我们主要从数据库的角度再来探讨同样的问题，即如何存储输入的数据，并在收到查询请求时，怎样重新找到数据。

作为一名应用系统开发人员，为什么要关注数据库内部的存储和检索呢？首先，你不太可能从头开始实现一套自己的存储引擎，往往需要从众多现有的存储引擎中选择一个适合自己应用的存储引擎。因此，为了针对你特定的工作负载而对数据库调优时，最好对存储引擎的底层机制有一个大概的了解。

特别地，针对事务型工作负载和针对分析型负载的存储引擎优化存在很大的差异。本章“事务处理与分析处理”和“面向列的存储”部分，将讨论一系列针对分析型进行优化的存储引擎。

我们首先讨论存储引擎，这些存储引擎用于大家比较熟悉的两种数据库，即传统的关系数据库和大多数所谓的NoSQL数据库。我们将研究两个存储引擎家族，即日志结构的存储引擎和面向页的存储引擎，比如B-tree。



# 数据库核心：数据结构

我们来看一个世界上最简单的数据库，它由两个Bash函数实现：

```
#!/bin/bash

db_set () {
    echo "$1,$2" >> database
}

db_get () {
    grep "^$1," database | sed -e "s/^$1,//" | tail -n 1
}
```

这两个函数实现了一个key-value存储。当调用`db_set key value`，它将在数据库中保存你所输入的key和value。key和value几乎可以是任何内容，例如，值可以是一个JSON文档。然后，调用`db_get key`，它会查找与输入key相关联的最新值并返回。

例如：

```
$ db_set 123456 '{"name":"London","attractions":["Big Ben","London Eye"]}'

$ db_set 42 '{"name":"San Francisco","attractions":["Golden Gate Bridge"]}'

$ db_get 42
{"name":"San Francisco","attractions":["Golden Gate Bridge"]}
```

它底层的存储格式其实非常简单：一个纯文本文件。其中每行包含一个key-value对，用逗号分隔（大致像一个CSV文件，忽略转义问题）。每次调用`db_set`即追加新内容到文件末尾，因此，如果多次更新某个键，旧版本的值不会被覆盖，而是需要查看文件中最后一次出现的键来找到最新的值（因此在`db_get`中使用`tail -n 1`）。

```
$ db_set 42 '{"name":"San Francisco","attractions":["Exploratorium"]}'

$ db_get 42
{"name":"San Francisco","attractions":["Exploratorium"]}

$ cat database
123456,{"name":"London","attractions":["Big Ben","London Eye"]}
42,{"name":"San Francisco","attractions":["Golden Gate Bridge"]}
42,{"name":"San Francisco","attractions":["Exploratorium"]}
```

对于简单的情况，追加到文件尾部方式通常足够高效，因而`db_set`函数性能很好。与`db_set`相似，许多数据库内部都使用日志（log），日志是一个仅支持追加式更新的数据文件。虽然真正的数据库有很多更为复杂问题需要考虑（例如并发控制、回收磁盘空间以控制日志文件大小、处理错误和部分完成写记录等），但是基本的原理是相同的。日志机制非常有用，本书的其余部分还将多次提到它。





日志这个词通常指的是应用程序的运行输出日志，来记录发生了什么事情。在本书中，日志则是一个更为通用的含义，表示一个仅能追加的记录序列集合。它可能是人类不可读的，可能是二进制格式而只能被其他程序来读取。

另一方面，如果日志文件保存了大量的记录，那么`db_get`函数的性能会非常差。每次想查找一个键，`db_get`必须从头到尾扫描整个数据库文件来查找键的出现位置。在算法术语中，查找的开销是 $O(n)$ ，即如果数据库的记录条数加倍，则查找需要两倍的时间。这一点并不好。

为了高效地查找数据库中特定键的值，需要新的数据结构：索引。在本章中，我们将介绍一些索引结构并对它们进行比较；它们背后的基本想法都是保留一些额外的元数据，这些元数据作为路标，帮助定位想要的数据库。如果希望用几种不同的方式搜索相同的数据，在数据的不同部分，我们可能定义多种不同的索引。

索引是基于原始数据派生而来的额外数据结构。很多数据库允许单独添加和删除索引，而不影响数据库的内容，它只会影响查询性能。维护额外的结构势必会引入开销，特别是在新数据写入时。对于写入，它很难超过简单地追加文件方式的性能，因为那已经是最简单的写操作了。由于每次写数据时，需要更新索引，因此任何类型的索引通常都会降低写的速度。

这里涉及存储系统中重要的权衡设计：适当的索引可以加速读取查询，但每个索引都会减慢写速度。为此，默认情况下，数据库通常不会对所有内容进行索引，它需要应用开发人员或数据库管理员，基于对应用程序典型查询模式的了解，来手动选择索引。目的是为应用程序提供最有利加速的同时，避免引入过多不必要的开销。

## 哈希索引

首先我们以键-值数据的索引开始。`key-value`类型并不是唯一可以索引的数据，但它随处可见，而且是其他更复杂索引的基础构造模块。

`key-value`存储与大多数编程语言所内置的字典结构非常相似，通常采用`hash map`（或者`hash table`，哈希表）来实现。许多算法资料中都介绍过`hash map`[1,2]，所以这里不再详细介绍它们工作的细节。那么，既然已经有了内存数据结构的`hash map`，为什么不用它们在磁盘上直接索引数据呢？

假设数据存储全部采用追加式文件组成，如之前的例子所示。那么最简单的索引策略就是：保存内存中的`hash map`，把每个键一一映射到数据文件中特定的字节偏移量，这样就可以找到每个值的位置，如图3-1所示。每当在文件中追加新的`key-value`对



时，还要更新hash map来反映刚刚写入数据的偏移量（包括插入新的键和更新已有的键）。当查找某个值时，使用hash map来找到文件中的偏移量，即存储位置，然后读取其内容。

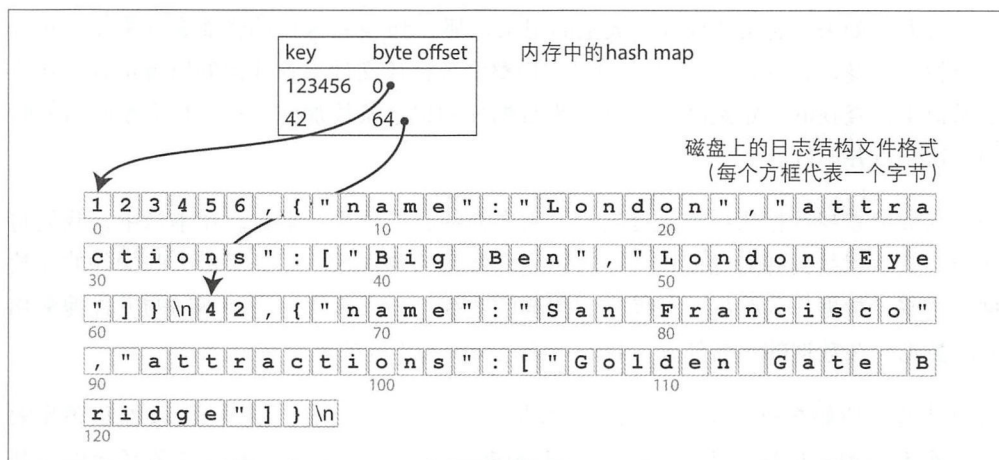


图3-1：采用类CSV格式来存储key-value对，以内存中的hash map来索引

这听起来可能过于简单，但它的确是一个可行的方法。事实上，这就是Bitcask（Riak中的默认存储引擎）所采用的核心做法<sup>[3]</sup>。Bitcask可以提供高性能的读和写，只要所有的key可以放入内存（因为hash map需要保存在内存中）。而value数据量则可以超过内存大小，只需一次磁盘寻址，就可以将value从磁盘加载到内存。如果那部分数据文件已经在文件系统的缓存中，则读取根本不需要任何的磁盘I/O。

像Bitcask这样的存储引擎非常适合每个键的值频繁更新的场景。例如，key可能是某个关于猫的视频URL，value是它播放的次数（每次有人单击播放按钮时就增加）。对于这种工作负载，有很多写操作，但是没有太多不同的key，即每个key都有大量的写操作，但将所有key保存在内存中是可行的。

如上所述，只追加到一个文件，那么如何避免最终用尽磁盘空间？一个好的解决方案是将日志分解成一定大小的段，当文件达到一定大小时就关闭它，并将后续写入到新的段文件中。然后可以在这些段上执行压缩，如图3-2所示。压缩意味着在日志中丢弃重复的键，并且只保留每个键最近的更新。

此外，由于压缩往往使得段更小（假设键在段内被覆盖多次），也可以在执行压缩的同时将多个段合并在一起，如图3-3所示。由于段在写入后不会进行修改，所以合并的段会被写入另一个新的文件。对于这些冻结段的合并和压缩过程可以在后台线程

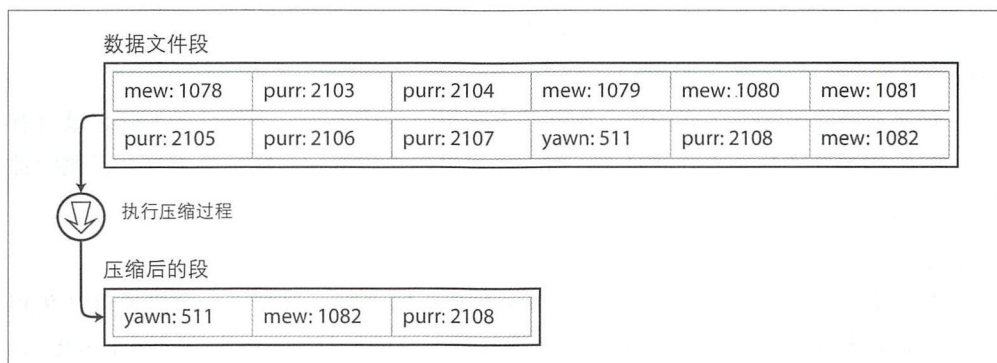


图3-2：压缩key-value更新日志文件（计算每个猫视频播放的次数），仅保留每个键的最新值中完成，而且运行时，仍然可以用旧的段文件继续正常读取和写请求。当合并过程完成后，将读取请求切换到新的合并段上，而旧的段文件可以安全删除。

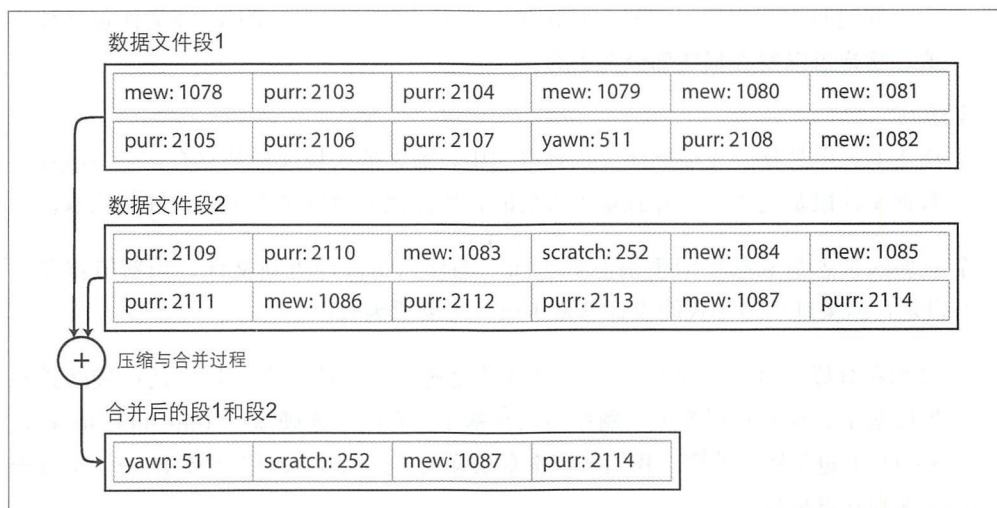


图3-3：同时执行段压缩和多个段合并

每个段现在都有自己的内存哈希表，将键映射到文件的偏移量。为了找到键的值，首先检查最新的段的hash map；如果键不存在，检查第二最新的段，以此类推。由于合并过程可以维持较少的段数量，因此查找通常不需要检查很多hash map。

还有很多细节方面的考虑才能使得这个简单的想法在实际中行之有效。简而言之，在真正地实现中有以下重要问题：

#### 文件格式

CSV不是日志的最佳格式。更快更简单的方法是使用二进制格式，首先以字节为





单位来记录字符串的长度，之后跟上原始字符串（不需要转义）。

### 删除记录

如果要删除键和它关联的值，则必须在数据文件中追加一个特殊的删除记录（有时候称为墓碑）。当合并日志段时，一旦发现墓碑标记，则会丢弃这个已删除键的所有值。

### 崩溃恢复

如果数据库重新启动，则内存中的hash map将丢失。原则上，可以通过从头到尾读取整个段文件，然后记录每个键的最新值的偏移量，来恢复每个段的hash map。但是，如果分段文件很大，可能扫描需要很长时间，这将使服务器重启变得缓慢。Bitcask通过将每个段的hash map的快照存储在磁盘上，可以更快地加载到内存中，以此加快恢复速度。

### 部分写入的记录

数据库随时可能崩溃，包括将记录追加到日志的过程中。Bitcask文件包括校验值，这样可以发现损坏部分并丢弃。

### 并发控制

由于写入以严格的先后顺序追加到日志中，通常的实现选择是只有一个写线程。数据文件段是追加的，并且是不可变的，所以他们可以被多个线程同时读取。

一个追加的日志乍看起来似乎很浪费空间：为什么不原地更新文件，用新值覆盖旧值？但是，结果证明追加式的设计非常不错，主要原因有以下几个：

- 追加和分段合并主要是顺序写，它通常比随机写入快得多，特别是在旋转式磁性硬盘上。在某种程度上，顺序写入在基于闪存的固态硬盘（solid state drives, SSD）上也是适合的<sup>[4]</sup>。我们将在本章后面的“比较B-tree和LSM-Trees”部分进一步讨论此问题。
- 如果段文件是追加的或不可变的，则并发和崩溃恢复要简单得多。例如，不必担心在重写值时发生崩溃的情况，留下一个包含部分旧值和部分新值混杂在一起的文件。
- 合并旧段可以避免随着时间的推移数据文件出现碎片化的问题。

但是，哈希表索引也有其局限性：

- 哈希表必须全部放入内存，所以如果有大量的键，就没那么幸运了。原则上，可以在磁盘上维护hash map，但不幸的是，很难使磁盘上的hash map表现良好。它



需要大量的随机访问I/O，当哈希变满时，继续增长代价昂贵，并且哈希冲突时需要复杂的处理逻辑<sup>[5]</sup>。

- 区间查询效率不高。例如，不能简单地支持扫描kitty00000和kitty99999区间内的所有键，只能采用逐个查找的方式查询每一个键。

在下一节中，我们将看到摆脱这些限制的其他索引结构。

## SSTables和LSM-Tree

在图3-3中，每个日志结构的存储段都是一组key-value对的序列。这些key-value对按照它们的写入顺序排列，并且对于出现在日志中的同一个键，后出现的值优于之前的值。除此之外，文件中key-value对的顺序并不重要。

现在简单地改变段文件的格式：要求key-value对的顺序按键排序。乍一看，这个要求似乎打破了顺序写规则，我们稍后会解释。

这种格式称为排序字符串表，或简称为*SSTable*。它要求每个键在每个合并的段文件中只能出现一次（压缩过程已经确保了）。SSTable相比哈希索引的日志段，具有以下优点：

1. 合并段更加简单高效，即使文件大于可用内存。方法类似于合并排序算法中使用的方法，如图3-4所示。并发读取多个输入段文件，比较每个文件的第一个键，把最小的键（根据排序顺序）拷贝到输出文件，并重复这个过程。这会产生一个新的按键排序的合并段文件。

如果相同的键出现在多个输入段怎么办？请记住，每个段包含在某段时间内写入数据库的所有值。这意味着一个输入段中的所有值肯定比其他段中的所有值更新（假设总是合并相邻的段）。当多个段包含相同的键时，可以保留最新段的值，并丢弃旧段中的值。

2. 在文件中查找特定的键时，不再需要在内存中保存所有键的索引。以图3-5为例，假设正在查找键handiwork，且不知道该键在段文件中的确切偏移。但是，如果知道键handbag和键handsome的偏移量，考虑到根据键排序，则键handiwork一定位于它们两者之间。这意味着可以跳到handbag的偏移，从那里开始扫描，直到找到handiwork（如果键handiwork不存在文件中，那么找不到）。

所以，仍然需要一个内存索引来记录某些键的偏移，但它可以是稀疏的，由于可



以很快扫描几千字节，对于段文件中每几千字节，只需要一个键就足够了<sup>注1</sup>。



图3-4：合并多个SSTable段，仅保留每个键的最新值

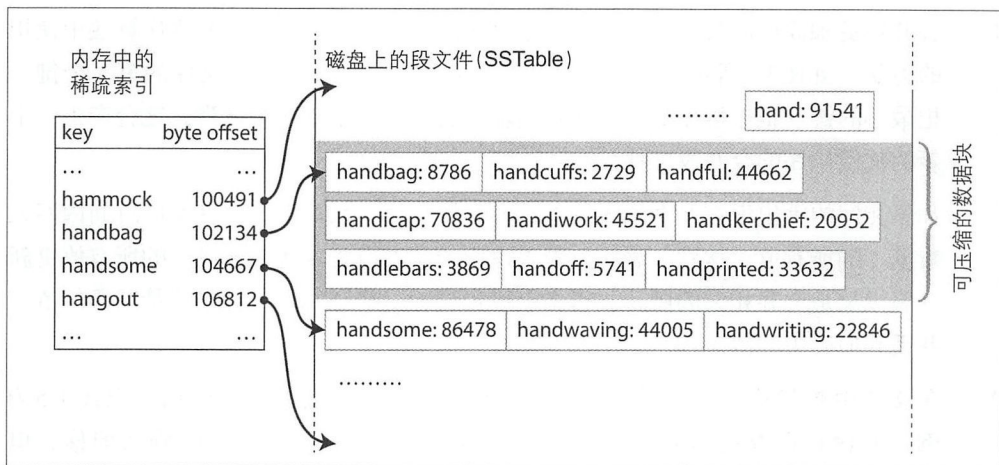


图3-5：SSTable及其内存中的索引

3. 由于读请求往往需要扫描请求范围内的多个key-value对，可以考虑将这些记录

注1： 如果所有键和值都具有固定的大小，则可以在分段文件上使用二分查找，并完全避免内存中的索引。但是，实践中通常是可变长度的，如果没有索引的话，难以确定一个记录的结束位置和下一个记录开始位置。





保存到一个块中并在写磁盘之前将其压缩（如图3-5中阴影区域所示）。然后稀疏内存索引的每个条目指向压缩块的开头。除了节省磁盘空间，压缩还减少了I/O带宽的占用。

## 构建和维护SSTables

到目前为止，似乎一切还算顺利，但是，考虑到写入可能以任意顺序出现，首先该如何让数据按键排序呢？

在磁盘上维护排序结构是可行的（参阅本章后面的“B-trees”），不过，将其保存在内存中更容易。内存排序有很多广为人知的树状数据结构，例如红黑树或AVL树<sup>[2]</sup>。使用这些数据结构，可以按任意顺序插入键并以排序后的顺序读取它们。

存储引擎的基本工作流程如下：

- 当写入时，将其添加到内存中的平衡树数据结构中（例如红黑树）。这个内存中的树有时被称为内存表。
- 当内存表大于某个阈值（通常为几兆字节）时，将其作为SSTable文件写入磁盘。由于树已经维护了按键排序的key-value对，写磁盘可以比较高效。新的SSTable文件成为数据库的最新部分。当SSTable写磁盘的同时，写入可以继续添加到一个新的内存表实例。
- 为了处理读请求，首先尝试在内存表中查找键，然后是最新的磁盘段文件，接下来是次新的磁盘段文件，以此类推，直到找到目标（或为空）。
- 后台进程周期性地执行段合并与压缩过程，以合并多个段文件，并丢弃那些已被覆盖或删除的值。

上述方案可以很好地工作。但它还存在一个问题：如果数据库崩溃，最近的写入（在内存表中但尚未写入磁盘）将会丢失。为了避免该问题，可以在磁盘上保留单独的日志，每个写入都会立即追加到该日志，就像上一节所述。那个日志文件不需要按键排序，这并不重要，因为它的唯一目的是在崩溃后恢复内存表。每当将内存表写入SSTable时，相应的日志可以被丢弃。

## 从SSTables到LSM-Tree

以上描述的算法本质上正是LevelDB<sup>[6]</sup>和RocksDB<sup>[7]</sup>所使用的，主要用于嵌入到其他应用程序的key-value存储引擎库。此外，在Riak中LevelDB可以用作Bitcask的替代品。类似的存储引擎还被用于Cassandra和HBase<sup>[8]</sup>，这两个引擎都受到Google的Bigtable论文<sup>[9]</sup>的启发（它引入了SSTable和内存表这两个术语）。



最初这个索引结构由Patrick O'Neil等人以日志结构的合并树（Log-Structured Merge-Tree，或LSM-Tree）<sup>[10]</sup>命名，它建立在更早期的日志结构文件系统之上<sup>[11]</sup>。因此，基于合并和压缩排序文件原理的存储引擎通常都被称为LSM存储引擎。

Lucene是Elasticsearch和Solr等全文搜索系统所使用的索引引擎，它采用了类似的方法来保存其词典<sup>[12,13]</sup>。全文索引比key-value索引复杂得多，但它基于类似的想法：给定搜索查询中的某个单词，找到提及该单词的所有文档（网页、产品描述等）。它主要采用key-value结构实现，其中键是单词（词条），值是所有包含该单词的文档ID的列表（倒排表）。在Lucene中，从词条到posting list的映射关系保存在类SSTable的排序文件中，这些文件可以根据需要在后台合并<sup>[14]</sup>。

## 性能优化

总是有很多细节值得深入优化，这样才能使存储引擎在实际中表现得更好。例如，当查找数据库中某个不存在的键时，LSM-Tree算法可能很慢：在确定键不存在之前，必须先检查内存表，然后将段一直回溯访问到最旧的段文件（可能必须从磁盘多次读取）。为了优化这种访问，存储引擎通常使用额外的布隆过滤器<sup>[15]</sup>（布隆过滤器是内存高效的数据结构，用于近似计算集合的内容。如果数据库中不存在某个键，它能够很快告诉你结果，从而节省了很多对于不存在的键的不必要的磁盘读取）。

还有不同的策略会影响甚至决定SSTables压缩和合并时的具体顺序和时机。最常见的方式是大小分级和分层压缩。LevelDB和RocksDB使用分层压缩（因此名称为LevelDB），HBase使用大小分级，Cassandra则同时支持这两种压缩<sup>[16]</sup>。在大小分级的压缩中，较新的和较小的SSTables被连续合并到较旧和较大的SSTables。在分层压缩中，键的范围分裂成多个更小的SSTables，旧数据被移动到单独的“层级”，这样压缩可以逐步进行并节省磁盘空间。

即使有许多细微的差异，但LSM-tree的基本思想（保存在后台合并的一系列SSTable）却足够简单有效。即使数据集远远大于可用内存，它仍然能够正常工作。由于数据按排序存储，因此可以有效地执行区间查询（从最小值到最大值扫描所有的键），并且由于磁盘是顺序写入的，所以LSM-tree可以支持非常高的写入吞吐量。

## B-trees

目前讨论的日志结构索引正在逐渐受到更多的认可，但它们还不是最常见的索引类型。最广泛使用的索引结构是另一种不同的：B-tree。

B-tree始见于1970年<sup>[17]</sup>，不到十年便被冠以“普遍存在”<sup>[18]</sup>，B-tree经受了长久的时



间考验。时至今日，它仍然是几乎所有关系数据库中的标准索引实现，许多非关系型数据库也经常使用。

像SSTable一样，B-tree保留按键排序的key-value对，这样可以实现高效的key-value查找和区间查询。但相似仅此而已：B-tree本质上具有非常不同的设计理念。

之前看到的日志结构索引将数据库分解为可变大小的段，通常大小为几兆字节或更大，并且始终按顺序写入段。相比之下，B-tree将数据库分解成固定大小的块或页，传统上大小为4 KB（有时更大），页是内部读/写的最小单元。这种设计更接近底层硬件，因为磁盘也是以固定大小的块排列。

每个页面都可以使用地址或位置进行标识，这样可以让一个页面引用另一个页面，类似指针，不过是指向磁盘地址，而不是内存。可以使用这些页面引用来构造一个树状页面，如图3-6所示。

某一页被指定为B-tree的根；每当查找索引中的一个键时，总是从这里开始。该页面包含若干个键和对子页的引用。每个孩子都负责一个连续范围内的键，相邻引用之间的键可以指示这些范围之间的边界。

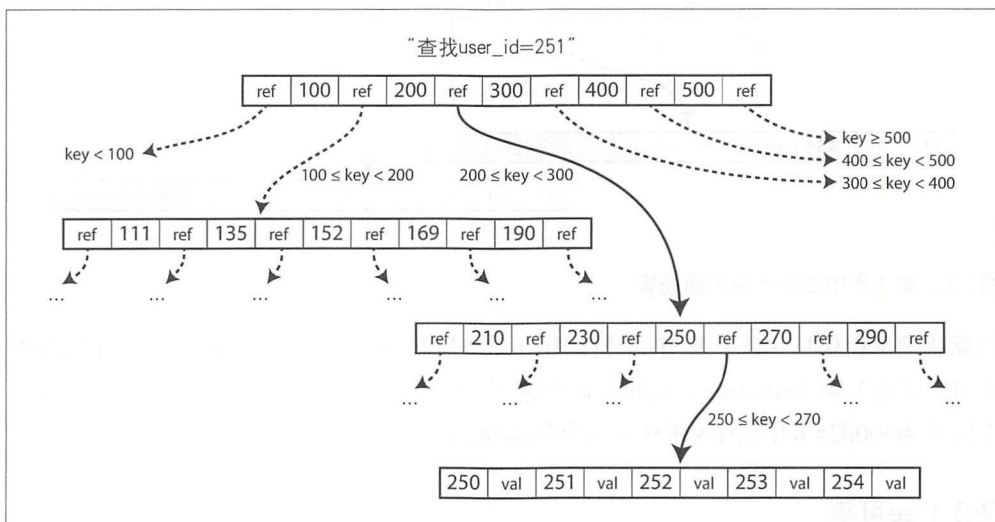


图3-6：使用B-tree索引查找关键字

在图3-6的例子中，假定正在查找键251，因此需要沿着200~300间的页引用，到达类似的页，它进一步将200~300范围分解成子范围。最终，我们到达一个包含单个键的页（叶子页），该页包含每个内联键的值或包含可以找到值的页的引用。





B-tree中一个页所包含的子页引用数量称为分支因子。例如，在图3-6中，分支因子为6。在实际中，分支因素取决于存储页面引用和范围边界所需的空间总量，通常为几百个。

如果要更新B-tree中现有键的值，首先搜索包含该键的叶子页，更改该页的值，并将页写回到磁盘（对该页的任何引用仍然有效）。如果要添加新键，则需要找到其范围包含新键的页，并将其添加到该页。如果页中没有足够的可用空间来容纳新键，则将其分裂为两个半满的页，并且父页也需要更新以包含分裂之后的新的键范围，如图3-7所示<sup>注2</sup>。

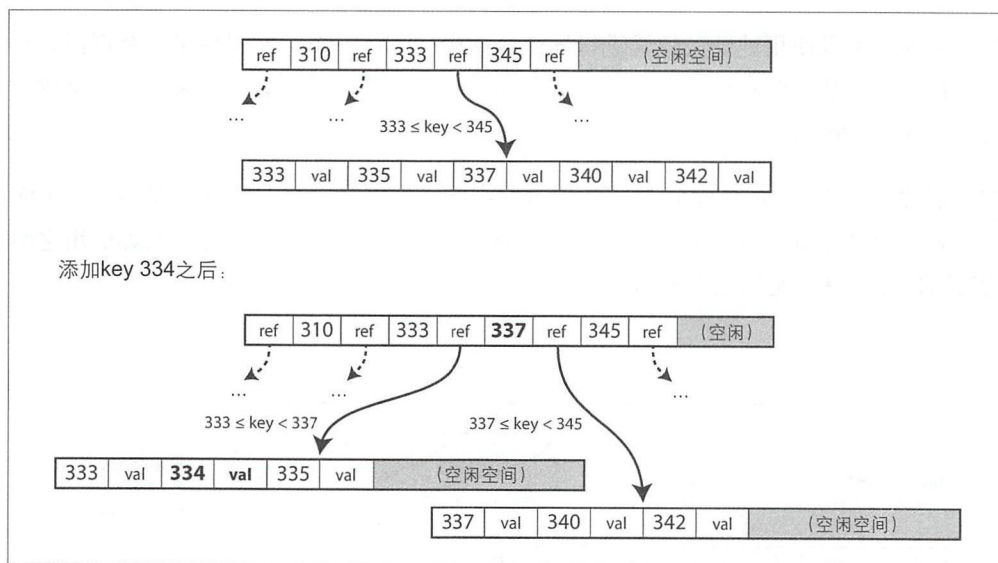


图3-7：插入B-tree时分裂页面场景

该算法确保树保持平衡：具有 $n$ 个键的B-tree总是具有 $O(\log n)$ 的深度。大多数数据库可以适合3~4层的B-tree，因此不需要遍历非常深的页面层次即可找到所需的页（分支因子为500的4 KB页的四级树可以存储高达256 TB）。

## 使B-tree可靠

B-tree底层的基本写操作是使用新数据覆盖磁盘上的旧页。它假设覆盖不会改变页的磁盘存储位置，也就是说，当页被覆盖时，对该页的所有引用保持不变。这与日志结构索引（如LSM-tree）形成鲜明对比，LSM-tree仅追加更新文件（并最终删除过时的文件），但不会修改文件。

注2：在B-tree中插入新的键比较直观，但是删除键（同时保持树平衡）有点复杂<sup>[2]</sup>。



可以认为磁盘上的页覆盖写对应确定的硬件操作。在磁性硬盘驱动器上，这意味着将磁头首先移动到正确的位置，然后旋转盘面，最后用新的数据覆盖相应的扇区。对于SSD，由于SSD必须一次擦除并重写非常大的存储芯片块，情况会更为复杂。

此外，某些操作需要覆盖多个不同的页。例如，如果插入导致页溢出，因而需分裂页，那么需要写两个分裂的页，并且覆盖其父页以更新对两个子页的引用。这是个比较危险的操作，因为如果数据库在完成部分页写入之后发生崩溃，最终会导致索引破坏（例如，可能有一个孤儿页，没有被任何其他页所指向）。

为了使数据库能从崩溃中恢复，常见B-tree的实现需要支持磁盘上的额外的数据结构：预写日志（write-ahead log, WAL），也称为重做日志。这是一个仅支持追加修改的文件，每个B-tree的修改必须先更新WAL然后再修改树本身的页。当数据库在崩溃后需要恢复时，该日志用于将B-tree恢复到最近一致的状态<sup>[5,20]</sup>。

原地更新页的另一个复杂因素是，如果多个线程要同时访问B-tree，则需要注意并发控制，否则线程可能会看到树处于不一致的状态。通常使用锁存器（轻量级的锁）保护树的数据结构来完成。在这方面，日志结构化的方法显得更简单，因为它们在后台执行所有合并，而不会干扰前端的查询，并且会不时地用新段原子地替换旧段。

## 优化B-tree

由于B-tree已经存在了很长时间，自然多年来开发了许多优化措施。这里只列举一些：

- 一些数据库（如LMDB）不使用覆盖页和维护WAL来进行崩溃恢复，而是使用写时复制方案<sup>[21]</sup>。修改的页被写入不同的位置，树中父页的新版本被创建，并指向新的位置。这种方法对于并发控制也很有帮助，详见后面第7章“快照隔离与可重复读”。
- 保存键的缩略信息，而不是完整的键，这样可以节省页空间。特别是在树中间的页中，只需要提供足够的信息来描述键的起止范围。这样可以将更多的键压入到页中，让树具有更高的分支因子，从而减少层数<sup>注3</sup>。
- 一般来说，页可以放在磁盘上的任何位置；没有要求相邻的页需要放在磁盘的相邻位置。如果查询需要按照顺序扫描大段的键范围，考虑到每个读取的页都可能需要磁盘I/O，所以逐页的布局可能是低效的。因此，许多B-tree的实现尝试对树

---

注3： 这个变种有时被称为B+树，然而优化是如此常见，以至于它通常不能与其他B-tree变种区分开来。



进行布局，以便相邻叶子页可以按顺序保存在磁盘上。然而，随着树的生长，维持这个顺序会变得越来越困难。相比之下，由于LSM-tree在合并过程中一次重写大量存储段，因此它们更容易让连续的键在磁盘上相互靠近。

- 添加额外的指针到树中。例如，每个叶子页面可能会向左和向右引用其同级的兄弟页，这样可以顺序扫描键，而不用跳回到父页。
- B-tree的变体如分形树<sup>[22]</sup>，借鉴了一些日志结构的想法来减少磁盘寻道（与分形无关）。

## 对比B-tree和LSM-tree

尽管B-tree的实现比LSM-tree的实现更为成熟，然而由于LSM-tree的性能特点，LSM-tree目前很有吸引力。根据经验，LSM-tree通常对于写入更快，而B-tree被认为对于读取更快<sup>[23]</sup>。读取通常在LSM-tree上较慢，因为它们必须在不同的压缩阶段检查多个不同的数据结构和SSTable。

然而，基准测试通常并不太确定，而且取决于很多工作负载的具体细节。最好测试特定工作负载，这样方便进行更有效的比较。在本节中，将简要讨论在测量存储引擎性能时值得考虑的一些要点。

### LSM-tree的优点

B-tree索引必须至少写两次数据：一次写入预写日志，一次写入树的页本身（还可能发生页分裂）。即使该页中只有几个字节更改，也必须承受写整个页的开销。一些存储引擎甚至覆盖相同的页两次，以避免在电源故障的情况下最终出现部分更新的页<sup>[24,25]</sup>。

由于反复压缩和SSTable的合并，日志结构索引也会重写数据多次。这种影响（在数据库内，由于一次数据库写入请求导致的多次磁盘写）称为写放大。对于SSD，由于只能承受有限次地擦除覆盖，因此尤为关注写放大指标。

对于大量写密集的应用程序，性能瓶颈很可能在于数据库写入磁盘的速率。在这种情况下，写放大具有直接的性能成本：存储引擎写入磁盘的次数越多，可用磁盘带宽中每秒可以处理的写入越少。

此外，LSM-tree通常能够承受比B-tree更高的写入吞吐量，部分是因为它们有时具有较低的写放大（尽管这取决于存储引擎的配置和工作负载），部分原因是它们以顺序方式写入紧凑的SSTable文件，而不必重写树中的多个页<sup>[26]</sup>。这种差异对于磁盘驱动器尤为重要，原因是磁盘的顺序写比随机写要快得多。





LSM-tree可以支持更好地压缩，因此通常磁盘上的文件比B-tree小很多。由于碎片，B-tree存储引擎使某些磁盘空间无法使用：当页被分裂或当一行的内容不能适合现有页时，页中的某些空间无法使用。由于LSM-tree不是面向页的，并且定期重写SSTables以消除碎片化，所以它们具有较低的存储开销，特别是在使用分层压缩时<sup>[27]</sup>。

在许多SSD上，固件内部使用日志结构化算法将随机写入转换为底层存储芯片上的顺序写入，所以存储引擎写入模式的影响不那么明显<sup>[19]</sup>。然而，更低的写放大和碎片减少对于SSD上仍然有益，以更紧凑的方式表示数据，从而在可用的I/O带宽中支持更多的读写请求。

## LSM-tree的缺点

日志结构存储的缺点是压缩过程有时会干扰正在进行的读写操作。即使存储引擎尝试增量地执行压缩，并且不影响并发访问，但由于磁盘的并发资源有限，所以当磁盘执行昂贵的压缩操作时，很容易发生读写请求等待的情况。这对吞吐量和平均响应时间的影响通常很小，但是如果观察较高的百分位数（参看第1章“描述性能”），日志结构化存储引擎的查询响应时间有时会相当高，而B-tree的响应延迟则更具确定性<sup>[28]</sup>。

高写入吞吐量时，压缩的另一个问题就会冒出来：磁盘的有限写入带宽需要在初始写入（记录并刷新内存表到磁盘）和后台运行的压缩线程之间所共享。写入空数据库时，全部的磁盘带宽可用于初始写入，但数据库的数据量越大，压缩所需的磁盘带宽就越多。

如果写入吞吐量很高并且压缩没有仔细配置，那么就会发生压缩无法匹配新数据写入速率的情况。在这种情况下，磁盘上未合并段的数量不断增加，直到磁盘空间不足，由于它们需要检查更多的段文件，因此读取速度也会降低。通常，即使压缩不能跟上，基于SSTable的存储引擎也不会限制到来的写入速率，因此需要额外的监控措施来及时发现这种情况<sup>[29,30]</sup>。

B-tree的优点则是每个键都恰好唯一对应于索引中的某个位置，而日志结构的存储引擎可能在不同的段中具有相同键的多个副本。如果数据库希望提供强大的事务语义，这方面B-tree显得更具有吸引力：在许多关系数据库中，事务隔离是通过键范围上的锁来实现的，并且在B-tree索引中，这些锁可以直接定义到树中<sup>[5]</sup>。在第7章，我们将更详细地讨论这一点。

B-tree在数据库架构中已经根深蒂固，为许多工作负载提供了一贯良好的性能，所以不太可能在短期内会消失。对于新的数据存储，日志结构索引则越来越受欢迎。不存



在快速和简单的规则来确定哪种存储引擎更适合你的用例，因此，实地的测试总是需要的。

## 其他索引结构

到目前为止，只讨论了key-value索引，它们像关系模型中的主键（primary key）索引。主键唯一标识关系表中的一行，或文档数据库中的一个文档，或图形数据库中的一个顶点。数据库中的其他记录可以通过其主键（或ID）来引用该行/文档/顶点，该索引用于解析此类引用。

二级索引也很常见。在关系数据库中，可以使用CREATE INDEX命令在同一个表上创建多个二级索引，并且它们通常对于高效地执行联结操作至关重要。例如，在第2章的图2-1中，很可能在user\_id列上有一个二级索引，以便可以在每个表中找到属于同一用户的所有行。

二级索引可以容易地基于key-value索引来构建。主要区别在于它的键不是唯一的，即可能有许多行（文档，顶点）具有相同键。这可以通过两种方式解决：使索引中的每个值成为匹配行标识符的列表（像全文索引中的posting list），或者追加一些行标识符来使每个键变得唯一。无论哪种方式，B-tree和日志结构索引都可以用作二级索引。

### 在索引中存储值

索引中的键是查询搜索的对象，而值则可以是以下两类之一：它可能是上述的实际行（文档，顶点），也可以是对其他地方存储的行的引用。在后一种情况下，存储行的具体位置被称为堆文件，并且它不以特定的顺序存储数据（它可以是追加的，或者记录删掉的行以便用新数据在之后覆盖它们）。堆文件方法比较常见，这样当存在多个二级索引时，它可以避免复制数据，即每个索引只引用堆文件中的位置信息，实际数据仍保存在一个位置。

当更新值而不更改键时，堆文件方法会非常高效：只要新值的字节数不大于旧值，记录就可以直接覆盖。如果新值较大，则情况会更复杂，它可能需要移动数据以得到一个足够大空间的新位置。在这种情况下，所有索引都需要更新以指向记录的新的堆位置，或者在旧堆位置保留一个间接指针<sup>[5]</sup>。

在某些情况下，从索引到堆文件的额外跳转对于读取来说意味着太多的性能损失，因此可能希望将索引行直接存储在索引中。这被称为聚集索引。例如，在MySQL InnoDB存储引擎中，表的主键始终是聚集索引，二级索引引用主键（而不是堆文件



位置)<sup>[31]</sup>。在SQL Server中，可以为每个表指定一个聚集索引<sup>[32]</sup>。

聚集索引（在索引中直接保存行数据）和非聚集索引（仅存储索引中的数据的引用）之间有一种折中设计称为覆盖索引或包含列的索引，它在索引中保存一些表的列值<sup>[33]</sup>。它可以支持只通过索引即可回答某些简单查询（在这种情况下，称索引覆盖了查询）<sup>[32]</sup>。

与任何类型的数据冗余一样，聚集和覆盖索引可以加快读取速度，但是它们需要额外的存储，并且会增加写入的开销。此外，数据库还需要更多的工作来保证事务性，这样应用程序不会因为数据冗余而得到不一致的结果。

## 多列索引

迄今为止讨论的索引只将一个键映射到一个值。如果需要同时查询表的多个列（或文档中的多个字段），那么这是不够的。

最常见的多列索引类型称为级联索引，它通过将一列追加到另一列，将几个字段简单地组合成一个键（索引的定义指定字段连接的顺序）。这就像一个老式的纸质电话簿，它提供从（lastname, firstname）到电话号码的索引。由于排序，索引可用于查找具有特定last name的所有人，或所有具有特定lastname-firstname组合的人。但是，它无法查找具有特定first name的所有人。

多维索引是更普遍的一次查询多列的方法，这对地理空间数据尤为重要。例如，餐馆搜索网站可能有一个包含每个餐厅的纬度和经度的数据库。当用户在地图上查看餐馆时，网站需要搜索用户正在查看的矩形地图区域内的所有餐馆。这要求一个二维的范围查询，如下所示：

```
SELECT * FROM restaurants WHERE latitude > 51.4946 AND latitude < 51.5079  
                                AND longitude > -0.1162 AND longitude < -0.1004;
```

标准B-tree或LSM-tree索引无法高效地应对这种查询，它只能提供一个纬度范围内（但在任何经度）的所有餐馆，或者所有经度范围内的餐厅（在北极和南极之间的任何地方），但不能同时满足。

一种选择是使用空格填充曲线将二维位置转换为单个数字，然后使用常规的B-tree索引<sup>[34]</sup>。更常见的是使用专门的空间索引，如R树。例如，PostGIS使用PostgreSQL的广义搜索树索引<sup>[35]</sup>实现了地理空间索引作为R树。篇幅所限，这里无法详细描述R树，但有很多关于它们的参考文献。

一个有趣的想法是，多维索引不仅仅适用于地理位置。例如，在电子商务网站上，可





以使用颜色维度（红色，绿色，蓝色）上的三维索引来搜索特定颜色范围内的产品，或在天气观测数据库中，可以在（日期，温度）上创建二维索引，以便高效地搜索在2013年中，温度在25~30°C之间的所有观测值。使用一维索引，将不得不从2013年扫描所有记录（无论温度如何），然后按温度进行过滤，反之亦然。二维索引可以通过时间戳和温度同时缩小查询范围。HyperDex使用了这种技术<sup>[36]</sup>。

## 全文搜索和模糊索引

到目前为止讨论的所有索引都假定具有确切的数据，并允许查询键的确切值或排序的键的取值范围。它们不支持搜索类似的键，如拼写错误的单词。这种模糊查询需要不同的技术。

例如，全文搜索引擎通常支持对一个单词的所有同义词进行查询，并忽略单词语法上的变体，在同一文档中搜索彼此接近的单词的出现，并且支持多种依赖语言分析的其他高级功能。为了处理文档或查询中的拼写错误，Lucene能够在某个编辑距离内搜索文本（编辑距离为1表示已经添加、删除或替换了一个字母）。

如上节“从SSTable到LSM-tree”所述，Lucene对其词典使用类似SSTable的结构。此结构需要一个小的内存索引来告诉查询，为了找到一个键，需要排序文件中的哪个偏移量。在LevelDB中，这个内存中的索引是一些键的稀疏集合，但是在Lucene中，内存中的索引是键中的字符序列的有限状态自动机，类似字典树<sup>[38]</sup>。这个自动机可以转换成Levenshtein自动机，它支持在给定编辑距离内高效地搜索单词<sup>[39]</sup>。

其他模糊搜索技术则沿着文档分类和机器学习的方向发展。有关更多细节，请参阅相关信息检索教程<sup>[40]</sup>。

## 在内存中保存所有内容

本章迄今为止讨论的数据结构都是为了适应磁盘限制。与内存相比，磁盘更难以处理。使用磁盘和SSD，如果要获得良好的读写性能，需要精心地安排磁盘上的数据布局。然而这些工作是值得的，因为磁盘有两个显著的优点：数据保存持久化（如果电源关闭，内容不会丢失），并且每GB容量的成本比内存低很多。

随着内存变得更便宜，每GB成本被摊薄。而许多数据集不是那么大，可以将它们完全保留在内存中，或者分布在多台机器上。这推动了内存数据库的发展。

一些内存中的key-value存储（如Memcached），主要用于缓存，如果机器重启造成的数据丢失是可以接受的。但是其他内存数据库旨在实现持久性，例如可以通过用特殊



硬件（如电池供电的内存），或者通过将更改记录写入磁盘，或者将定期快照写入磁盘，以及复制内存中的状态到其他机器等方式来实现。

当内存数据库重启时，它需要重新载入其状态，无论是从磁盘还是通过网络从副本（除非使用特殊硬件）。尽管写入磁盘，但磁盘仅仅用作为了持久性目的的追加日志，读取完全靠内存服务。此外，写入磁盘还具有一些运维方面优势：磁盘上的文件可以容易地通过外部工具来执行备份、检查和分析。

诸如VoltDB、MemSQL和Oracle TimesTen的产品是具有关系模型的内存数据库，相关供应商声称通过移除所有与管理磁盘数据结构相关的开销，它们可以获得极大的性能提升<sup>[41, 42]</sup>。RAMCloud是一个开源的、具有持久性的内存key-value存储（对内存和磁盘上的数据使用日志结构）<sup>[43]</sup>。而Redis和Couchbase通过异步写入磁盘提供较弱的持久性。

与直觉相反，内存数据库的性能优势并不是因为它们不需要从磁盘读取。如果有足够的内存，即使是基于磁盘的存储引擎，也可能永远不需要从磁盘读取，因为操作系统将最近使用的磁盘块缓存在内存中。相反，内存数据库可以更快，是因为它们避免使用写磁盘的格式对内存数据结构编码的开销<sup>[44]</sup>。

除了性能外，内存数据库的另一个有意思的地方是，它提供了基于磁盘索引难以实现的某些数据模型。例如，Redis为各种数据结构（如优先级队列和集合）都提供了类似数据库的访问接口。由于所有的数据都保存在内存中，所以实现可以比较简单。

最近的研究表明，内存数据库架构可以扩展到支持远大于可用内存的数据集，而不会导致以磁盘为中心架构的开销<sup>[45]</sup>。所谓的反缓存方法，当没有足够的内存时，通过将最近最少使用的数据从内存写到磁盘，并在将来再次被访问时将其加载到内存。这与操作系统对虚拟内存和交换文件的操作类似，但数据库可以在记录级别而不是整个内存页的粒度工作，因而比操作系统更有效地管理内存。不过，这种方法仍然需要索引完全放入内存（如本章开头的Bitcask示例）。

如果将来非易失性存储（non-volatile memory, NVM）技术得到更广泛普及，可能还需要进一步改变存储引擎设计<sup>[46]</sup>。目前这是一个新的研究领域，但值得密切关注。

## 事务处理与分析处理

在商业数据处理的早期阶段，写入数据库通常对应于商业交易场景，例如销售、订单、支付员工工资等。尽管后来数据库扩展到了不涉及金钱交易的领域，事务一词仍然存在，主要指组成一个逻辑单元的一组读写操作。



事务不一定具有ACID（原子性、一致性、隔离性和持久性）属性。事务处理只是意味着允许客户端进行低延迟读取和写入，相比于只能周期性地运行（如每天一次）的批处理作业。我们将在第7章讨论ACID属性，在第10章讨论批处理。

尽管数据库开始被用于许多不同种类的数据，例如博客的评论、游戏中的动作、通讯录中的联系人等，然而其基本访问模式仍然与处理业务交易类似。应用程序通常使用索引中的某些键查找少量记录。根据用户的输入插入或更新记录。因为这些应用程序是交互式的，所以访问模式被称为在线事务处理（online transaction processing，OLTP）。

然而，数据库也开始越来越多地用于数据分析，数据分析具有非常不同的访问模式。通常，分析查询需要扫描大量记录，每个记录只读取少数几列，并计算汇总统计信息（如计数、求和或平均值），而不是返回原始数据给用户。例如，如果数据是销售交易表，那么分析查询可能包括：

- 一月份每个店铺的总收入是多少？
- 在最近促销的期间，比平时多卖了多少香蕉？
- 哪个品牌的婴儿食品最常与某品牌的尿布一起购买？

这些查询通常由业务分析师编写，以形成有助于公司管理层更好决策（商业智能）的报告。为了区分使用数据库与事务处理的模式，称之为在线分析处理（online analytic processing，OLAP）<sup>[47]注4</sup>。OLTP和OLAP之间的区别有时并不那么明确，它们的一些典型的特性总结见表3-1。

表3-1：对比事务处理与分析系统的主要特性

属性	事务处理系统（OLTP）	分析系统（OLAP）
主要读特征	基于键，每次查询返回少量的记录	对大量记录进行汇总
主要写特征	随机访问，低延迟写入用户的输入	批量导入（ETL）或事件流
典型使用场景	终端用户，通过网络应用程序	内部分析师，为决策提供支持
数据表征	最新的数据状态（当前时间点）	随着时间而变化的所有事件历史
数据规模	GB到TB	TB到PB

注4： OLAP中的在线含义尚不清楚。它可能指的是这样的事实，查询不仅仅针对预定义的报告，而且分析人员以交互方式使用OLAP系统进行探索性查询。





最初，相同的数据库可以同时用于事务处理和分析查询。在这方面，SQL被证明是非常灵活的，可以同时胜任OLTP类型和OLAP类型查询。然而，在20世纪80年代后期和90年代初期的一种趋势是，公司放弃使用OLTP系统用于分析目的，而是在单独的数据库上运行分析。这个单独的数据库被称为数据仓库。

## 数据仓库

企业可能有几十种不同的交易处理系统，例如面向客户的网站提供支持的系统、控制实体店的销售（结账）系统、跟踪仓库库存、为车辆规划路线、供应商管理、员工管理等。这些系统中的每一个都足够复杂，往往需要一个专门团队来维护，最终导致这些系统彼此独立运行。

由于这些OLTP系统对于业务的运行至关重要，所以往往期望它们高度可用，处理事务时延迟足够低，并且数据库管理员要密切关注OLTP数据库运行状态。数据库管理员通常不愿意让业务分析人员在OLTP数据库上直接运行临时分析查询，这些查询通常代价很高，要扫描大量数据集，这可能会损害并发执行事务的性能。

相比之下，数据仓库则是单独的数据库，分析人员可以在不影响OLTP操作的情况下尽情地使用<sup>[48]</sup>。数据仓库包含公司所有各种OLTP系统的只读副本。从OLTP数据库（使用周期性数据转储或连续更新流）中提取数据，转换为分析友好的模式，执行必要的清理，然后加载到数据仓库中。将数据导入数据仓库的过程称为提取-转换-加载（Extract-Transform-Load，ETL），如图3-8所示。

几乎所有的大型企业都有数据仓库，但是在小型企业中却几乎闻所未闻。这可能是因为大多数小公司没有那么多不同的OLTP系统，大多数小公司只拥有少量的数据，完全可以在传统的SQL数据库中直接进行查询分析，甚至可以在电子表格中进行分析。对于大公司，需要做大量繁重的工作来完成在小公司看似很简单的一些事情。

使用单独的数据仓库而不是直接查询OLTP系统进行分析，很大的优势在于数据仓库可以针对分析访问模式进行优化。事实证明，本章前半部分讨论的索引算法适合OLTP，但不擅长应对分析查询。

在本章的其余部分，将重点讨论针对分析型而优化的存储引擎。

### OLTP数据库和数据仓库之间的差异

数据仓库的数据模型最常见的是关系型，因为SQL通常适合分析查询。有许多图形化数据分析工具，它们可以生成SQL查询、可视化结果并支持分析师探索数据，例如通过诸如向下钻取、切片和切丁等操作。

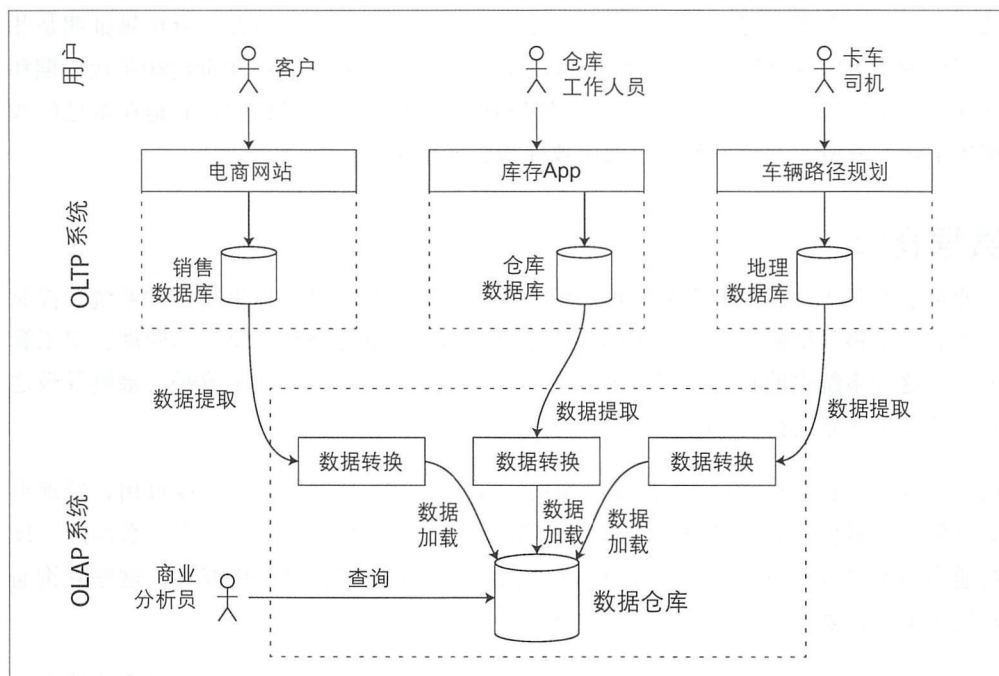


图3-8：数据仓库和简化的ETL过程

表面上，数据仓库和关系型OLTP数据库看起来相似，因为它们都具有SQL查询接口。然而，系统内部实则差异很大，它们针对迥然不同的查询模式进行了各自优化。许多数据库供应商现在专注于支持事务处理或分析工作负载，但不能同时支持两者。

一些数据库（如Microsoft SQL Server和SAP HANA）在同一产品中支持事务处理和数据仓库。然而，它们越来越成为两个独立的存储和查询引擎，这些引擎恰好可以通过一个通用的SQL界面进行访问[49-51]。

诸如Teradata、Vertica、SAP HANA和ParAccel等数据仓库供应商，通过昂贵的商业许可来销售其系统。Amazon RedShift是ParAccel的托管版本。最近，还出现了大量开源的基于Hadoop的SQL项目，它们还比较年轻，但都试图与商业数据仓库系统展开竞争。包括Apache Hive、Spark SQL、Cloudera Impala、Facebook Presto、Apache Tajo和Apache Drill<sup>[52,53]</sup>。其中一些系统是基于Google Dremel而构建的<sup>[54]</sup>。

## 星型与雪花型分析模式

如第2章所述，根据不同的应用需求，事务处理领域广泛使用了多种不同数据模型。而另一方面，分析型业务的数据模型则要少得多。许多数据仓库都相当公式化的使用



了星型模式，也称为维度建模<sup>[55]</sup>。

图3-9所示的模式可用于零售数据仓库。模式的中心是一个所谓的事实表（在这个例子中，它被称为fact\_sales）。事实表的每一行表示在特定时间发生的事件（这里，每一行代表客户购买的一个产品）。如果我们正在分析网站流量而不是零售，则每一行可能表示页面视图或用户的单击。

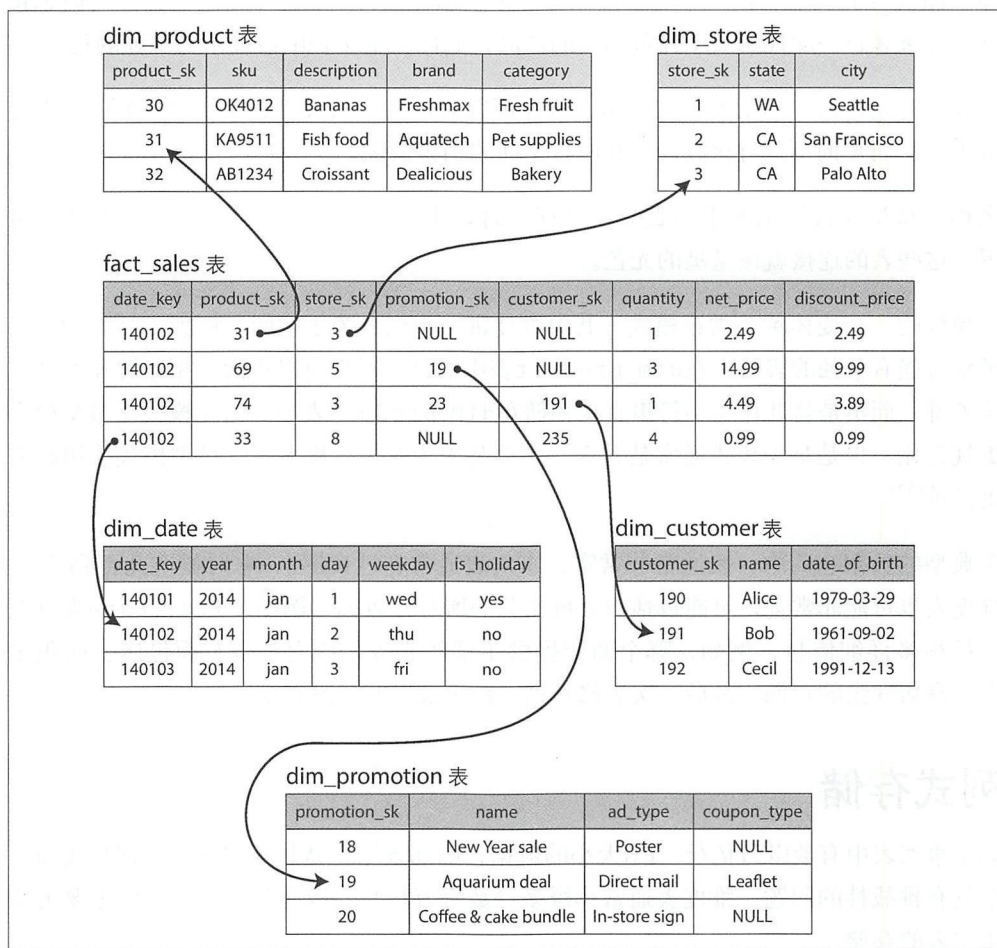


图3-9：用于数据仓库的星型模式示例

通常，事实被捕获为单独的事件，这样之后的分析具有最大的灵活性。不过，这意味着事实表可能会变得非常庞大。像苹果、沃尔玛或者eBay这样的大企业，其数据仓库可能有数十PB的交易历史，其中大部分都保存在事实表中<sup>[56]</sup>。

事实表中的列是属性，例如产品销售的价格和供应商处购买的成本（可以计算出利润





率)。其他列可能会引用其他表的外键，称为维度表。由于事实表中的每一行都代表一个事件，维度通常代表事件的对象(who)、什么(what)、地点(when)、时间(when)、方法(how)以及原因(why)。

例如，在图3-9中，其中一个维度是销售的产品。dim\_product表中的每一行代表一种出售的产品，包括库存单位(SKU)、说明、品牌名称、类别、脂肪含量、包装尺寸等。fact\_sales表中的每一行都使用外键来表示在该特定事务中出售的产品(简单起见，如果客户一次性购买了几种不同的产品，它们在事实表中被表示为单独的行)。

日期和时间通常使用维度表来表示，这样可以对日期(如公共假期)的相关信息编码，从而查询可以对比假期和非假日之间的销售情况。

名称“星型模式”来源于当表关系可视化时，事实表位于中间，被一系列维度表包围；这些表的连接就像星星的光芒。

该模板的一个变体称为雪花模式，其中维度进一步细分为子空间。例如，品牌和产品类别可能有单独的表格，在dim\_product表中的每一行都可以再次引用品牌和类别作为外键，而不是将其作为字符串直接存储在dim\_product表中。雪花模式比星型模式更规范化，但是星型模式通常是首选，主要是因为对于分析人员，星型模式使用起来更简单<sup>[55]</sup>。

在典型的数据仓库中，表通常非常宽：事实表通常超过100列，有时候有几百列<sup>[51]</sup>。维度表也可能非常宽，可能包括与分析相关的所有元数据，例如，dim\_store表可能包括很多详细信息，例如，哪个商店提供了哪些服务，店内是否有面包店、面积多大、商店开张的日期、最后一次装修日期、距离最近的公路有多远等。

## 列式存储

如果事实表中有数以万亿行、PB大小的数据，则高效地存储和查询这些数据将成为一个具有挑战性的问题。维度表通常小得多(数百万行)，因此在本节中，将主要关注事实表的存储。

虽然事实表通常超过100列，但典型的数据仓库查询往往一次只访问其中的4或5个("SELECT \*"查询很少用于分析)<sup>[51]</sup>。下面示例3-1中的查询中，它会访问大量行(某人在2013年中每一次购买水果或糖果)，但结果只需要返回fact\_sales表的三列：date\_key、product\_sk和quantity，其他列不在输出范围。



示例3-1：分析人们购买新鲜水果或糖果的倾向是否取决于一周中的某天

```
SELECT
    dim_date.weekday, dim_product.category,
    SUM(fact_sales.quantity) AS quantity_sold
FROM fact_sales
    JOIN dim_date ON fact_sales.date_key = dim_date.date_key
    JOIN dim_product ON fact_sales.product_sk = dim_product.product_sk
WHERE
    dim_date.year = 2013 AND
    dim_product.category IN ('Fresh fruit', 'Candy')
GROUP BY
    dim_date.weekday, dim_product.category;
```

如何高效地执行这个查询？

在大多数OLTP数据库中，存储以面向行的方式布局：来自表的一行的所有值彼此相邻存储。文档数据库也是类似，整个文档通常被存储为一个连续的字节序列，例如图3-1的CSV示例。

为了处理像示例3-1这样的查询，可以在`fact_sales.date_key`和/或`fact_sales.product_sk`上使用索引，告诉存储引擎在哪里查找特定日期或特定产品的所有销售。但是，面向行的存储引擎仍然需要将所有行（每个由超过100个属性组成）从磁盘加载到内存中、解析它们，并过滤出不符合所需条件的行。这可能需要很长时间。

面向列存储的想法很简单：不要将一行中的所有值存储在一起，而是将每列中的所有值存储在一起。如果每个列存储在一个单独的文件中，查询只需要读取和解析在该查询中使用的那些列，这可以节省大量的工作。该原理如图3-10所示。



列存储在关系数据模型中最容易理解，但它同样适用于非关系数据。例如，Parquet<sup>[57]</sup>是基于Google的Dremel<sup>[54]</sup>的一种支持文档数据模型的列存储格式。

面向列的存储布局依赖一组列文件，每个文件以相同顺序保存着数据行。因此，如果需要重新组装整行，可以从每个单独的列文件中获取第23个条目，并将它们放在一起构成表的第23行。

## 列压缩

除了仅从磁盘中加载查询所需的列之外，还可以通过压缩数据来进一步降低对磁盘吞吐量的要求。幸运的是，面向列的存储恰好非常适合压缩。



看看图3-10中每列的值序列：它们看起来有很多重复，这是压缩的好兆头。取决于列中具体数据模式，可以采用不同的压缩技术。在数据仓库中特别有效的一种技术是位图编码，如图3-11所示。

fact_sales 表							
date_key	product_sk	store_sk	promotion_sk	customer_sk	quantity	net_price	discount_price
140102	69	4	NULL	NULL	1	13.99	13.99
140102	69	5	19	NULL	3	14.99	9.99
140102	69	5	NULL	191	1	14.99	14.99
140102	74	3	23	202	5	0.99	0.89
140103	31	2	NULL	NULL	1	2.49	2.49
140103	31	3	NULL	NULL	3	14.99	9.99
140103	31	3	21	123	1	49.99	39.99
140103	31	8	NULL	233	1	0.99	0.99

列式存储磁盘布局：

date\_key 文件内容： 140102, 140102, 140102, 140102, 140103, 140103, 140103, 140103

product\_sk 文件内容： 69, 69, 69, 74, 31, 31, 31, 31

store\_sk 文件内容： 4, 5, 5, 3, 2, 3, 3, 8

promotion\_sk 文件内容： NULL, 19, NULL, 23, NULL, NULL, 21, NULL

customer\_sk 文件内容： NULL, NULL, 191, 202, NULL, NULL, 123, 233

quantity 文件内容： 1, 3, 1, 5, 1, 3, 1, 1

net\_price 文件内容： 13.99, 14.99, 14.99, 0.99, 2.49, 14.99, 49.99, 0.99

discount\_price 文件内容： 13.99, 9.99, 14.99, 0.89, 2.49, 9.99, 39.99, 0.99

图3-10：按列而不是按行存储关系数据

通常，列中的不同值的数量小于行数（一个例如，零售商可能拥有数十亿个销售交易，但只有100000个不同的产品）。现在可以使用 $n$ 个不同值的列，并将其转换为 $n$ 个单独的位图：一个位图对应每个不同的值，一个位对应一行。如果行具有该值，该位为1，否则为0。

如果 $n$ 非常小（例如，表示国家的列可能具有大约200个不同的值），那么这些位图由每行一位存储。但是，如果 $n$ 越大，在大多数位图中将会有很多零（它们很稀疏）。此时，位图也可以进行游程编码，如图3-11底部所示。这样列的编码非常紧凑。

这些位图索引非常适合在数据仓库中常见的查询。例如：

WHERE product\_sk IN (30, 68, 69):

加载product\_sk = 30、product\_sk = 68和product\_sk = 69的三个位图，并计算三个位图的按位或，这可以非常高效地完成。



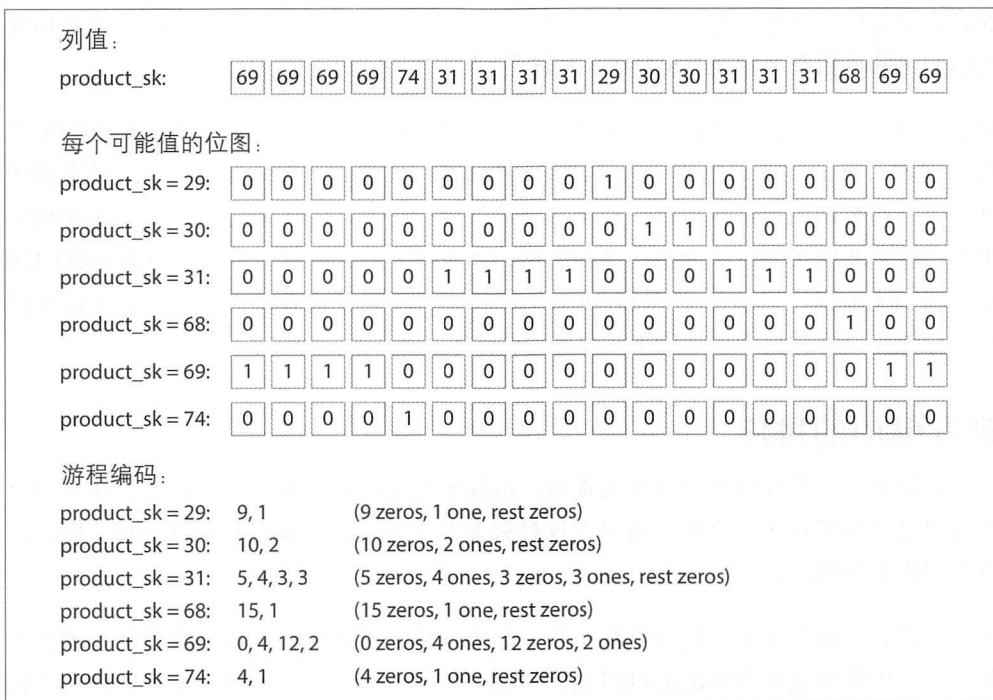


图3-11：压缩的位图索引存储单列

WHERE product\_sk = 31 AND store\_sk = 3:

加载product\_sk = 31和store\_sk = 3的位图，并按位与计算。这样做也是可行的，这是因为这些列包含相同顺序的行，因此一系列的位图中的第k位对应于另一列的位图中与第k位相同的行。

对于不同类型的数据，还有各种其他压缩方法<sup>[58]</sup>。



### 面向列的存储和列族

Cassandra和HBase有一个列族的概念，它们继承自Google Bigtable[9]。但是，将它们称为面向列则非常令人误解：在每个列族中，它们将一行中的所有列与行主键一起保存，并且不使用列压缩。因此，Bigtable模型仍然主要是面向行。

## 内存带宽和矢量化处理

对于需要扫描数百万行的数据仓库查询，将数据从磁盘加载到内存的带宽是一大瓶颈。然而，这还不是唯一的瓶颈。分析数据库的开发人员还要关心如何高效地将内存



的带宽用于CPU缓存，避免分支错误预测和CPU指令处理流水线中的气泡，并利用现代CPU中的单指令多数据（SIMD）指令<sup>[59, 60]</sup>。

除了减少需要从磁盘加载的数据量之外，面向列的存储布局也有利于高效利用CPU周期。例如，查询引擎可以将一大块压缩列数据放入CPU的L1缓存中，并以紧凑循环（即没有函数调用）进行迭代。对于每个被处理的记录，CPU能够比基于很多函数调用和条件判断的代码更快地执行这种循环。列压缩使得列中更多的行可以加载到L1缓存。诸如先前描述的按位AND和OR的运算符，可被设计成直接对这样的列压缩数据块进行操作。这种技术被称为矢量化处理<sup>[49, 58]</sup>。

## 列存储中的排序

在列存储中，行的存储顺序并不太重要。最简单的是按插入顺序保存，这样插入一个新行只是追加到每个列文件。也可以选择强制某个顺序，就像之前SSTable一样，并将其用作索引机制。

但请注意，单独排序每列是没有意义的，如果这样的话就无法知道列中的某一项属于哪一行。因为知道某列中的第 $k$ 项和另一列的第 $k$ 项一定属于同一行，基于这种约定我们可以重建一行。

相反，即使数据是按列存储的，它也需要一次排序整行。数据库管理员可以基于常见查询的知识来选择要排序表的列。例如，如果查询经常以日期范围为目标，例如上个月，那么明智的做法是将`date_key`设置为第一个排序键。查询优化器只扫描上个月的行为，这比起扫描所有行快得多。

当第一列排序出现相同值时，可以指定第二列继续进行排序。例如，如果`date_key`是图3-10中的第一个排序键，然后可以指定`product_sk`成为第二个排序键，这样同一天同一产品的所有销售在存储中被分组在一起。这有助于在某个日期范围内按产品进行分组或过滤的查询。

排序的另一个优点是它可以帮助进一步压缩列。如果主排序列上没有很多不同的值，那么在排序之后，它将出现一个非常长的序列，其中相同的值在一行中重复多次。一个简单的游程编码，如图3-11中的位图那样，即使该表可能拥有数十亿行，也可以将其压缩到几千字节。

基于第一个排序键的压缩效果通常最好。第二个和第三个排序键会使情况更加复杂，也通常不会有太多相邻的重复值。排序优先级进一步下降的列基本上会呈现接近随机的顺序，因此通常无法压缩。但总体来讲，对前几列排序仍然可以获得不错的收益。



## 几种不同的排序

C-Store最早引入了一种改进想法，并被商业数据仓库Vertica<sup>[61,62]</sup>所采用。考虑到不同的查询会从不同的排序中获益，那么为什么不以多种不同的方式存储相同的数据呢？无论如何，数据需要复制到多台机器，这样在一台机器发生故障时，不会丢失数据。不妨存储不同方式排序的冗余数据，以便在处理查询时，可以选择最适合特定查询模式的排序版本。

面向列的存储具有多个排序顺序，这有些类似在面向行的存储中具有多个二级索引。但最大的区别是，面向行的存储将每一行都保存在一个位置（在堆文件或聚集索引中），而二级索引只包含指向匹配行的指针。而对于列存储，通常没有任何指向别处数据的指针，只有包含值的列。

## 列存储的写操作

上述这些优化对于数据仓库很有意义，因为大多数负载由分析人员运行的大型只读查询组成。面向列的存储、压缩和排序都非常有助于加速读取查询。但是，它们的缺点是让写入更加困难。

像B-tree使用的原地更新方式，对于压缩的列是不可能的。如果在排序表的中间插入一行，那么很可能不得不重写所有的列文件。因为各行是由它们在列中的位置标识的，所以插入操作必须一致地更新所有列。

幸运的是，在本章前面已经看到了一个很好的解决方案LSM-tree。所有的写入首先进入内存存储区，将其添加到已排序的结构中，接着再准备写入磁盘。内存中的存储是面向行还是面向列无关紧要。当累积了足够的写入时，它们将与磁盘上的列文件合并，并批量写入新文件。本质上这正是Vertica所做的事情<sup>[62]</sup>。

执行查询时，需要检查磁盘上的列数据和内存中最近的写入，并结合这两者。而查询优化器可以对用户隐藏这些内部细节。从数据分析师的角度来看，插入、更新或删除数据可以立即反映在随后的查询中。

## 聚合：数据立方体与物化视图

未必每个数据仓库都基于列存储，传统的面向行数据库以及其他一些架构也有使用。然而，对于临时分析查询，列存储性能要快得多，因此它正在迅速普及<sup>[51,63]</sup>。

数据仓库的另一个值得一提的是物化聚合。如前所述，数据仓库查询通常涉及聚合函数，例如SQL中的COUNT、SUM、AVG、MIN或MAX。如果许多不同查询使用相同





的聚合，每次都处理原始数据将非常浪费。为什么不缓存查询最常使用的一些计数或总和呢？

创建这种缓存的一种方式物化视图。在关系数据模型中，它通常被定义为标准（虚拟）视图：一个类似表的对象，其内容是一些查询的结果。不同的是，物化视图是查询结果的实际副本，并被写到磁盘，而虚拟视图只是用于编写查询的快捷方式。从虚拟视图中读取时，SQL引擎将其动态地扩展到视图的底层查询，然后处理扩展查询。

当底层数据发生变化时，物化视图也需要随之更新，因为它是数据的非规范化副本。数据库可以自动执行，但这种更新方式会影响数据写入性能，这就是为什么在OLTP数据库中不经常使用物化视图的原因。而对于大量读密集的数据仓库，物化视图则更有意义（它们是否能够真正地提高性能还要取决于具体情况）。

物化视图常见的一种特殊情况称为数据立方体或OLAP立方体<sup>[64]</sup>。它是由不同维度分组的聚合网格，如图3-12所示的例子。

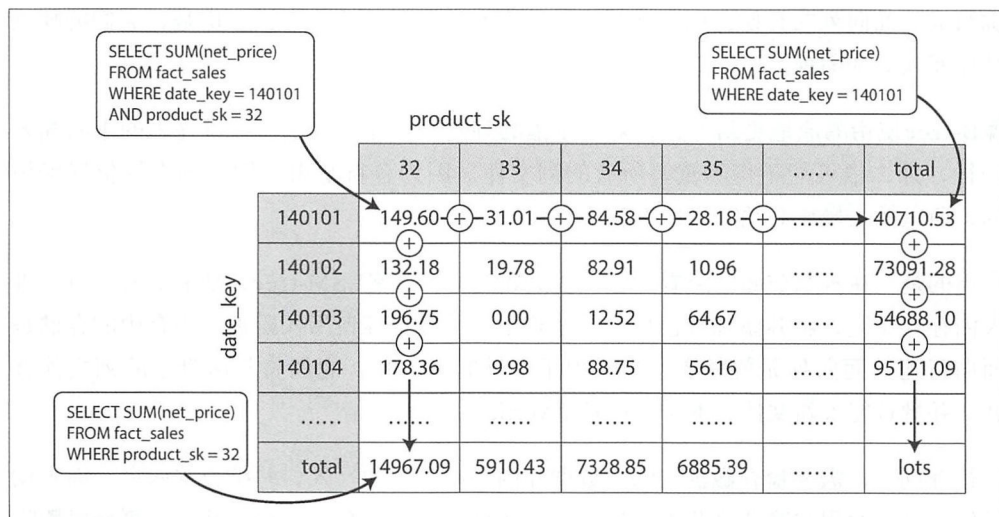


图3-12：数据立方体的两个维度，通过求和聚合数据

想象一下，每个事实只包含两个维度表的外键，在图3-12中，它们是日期（date\_key）和产品（product\_sk）。现在绘制这样一个二维表，日期沿着一个轴，产品沿着另一个轴。每个单元格即是date-product组合的所有事实的属性（例如，net\_price）的聚合（例如，SUM）。然后，可以沿着每一行或列应用聚合操作，得到一个减少一个维度的总和（按产品销售额而不管日期，或按日期的销售额而不管产品）。



一般来说，事实表的维度不止两个，例如在图3-9中有五个维度：日期、产品、商店、促销和客户。想象多维超立方体是什么样子有些困难，但是原理是类似的：每个单元格包含特定日期-产品-商店-促销-客户组合的销售值。然后可以沿着每个维度汇总这些值。

物化数据立方体的优点是某些查询会非常快，主要是它们已被预先计算出来。例如，如果想知道昨天每个商店的总销售量，只需要直接查看对应维度的总和，而不需要扫描数百万行。

缺点则是，数据立方体缺乏像查询原始数据那样的灵活性。例如，因为价格不是其中的一个维度，所以没有办法直接计算成本超过100美元的物品所占销售的比重。因此，大多数数据仓库都保留尽可能多的原始数据，仅当数据立方体可以对特定查询显著提升性能时，才会采用多维数据聚合。

## 小结

本章我们简单介绍了数据库内部如何处理存储与检索。诸如，数据库存储新数据时会发生什么，以及之后查询数据时，数据库会做什么？

概括来讲，存储引擎分为两大类：针对事务处理（OLTP）优化的架构，以及针对分析型（OLAP）的优化架构。它们典型的访问模式存在很大差异：

- OLTP系统通常面向用户，这意味着它们可能收到大量的请求。为了处理负载，应用程序通常在每个查询中只涉及少量的记录。应用程序基于某种键来请求记录，而存储引擎使用索引来查找所请求键的数据。磁盘寻道时间往往是瓶颈。
- 由于不是直接面对最终用户，数据仓库和类似的分析型系统相对并不太广为人知，它们主要由业务分析师使用。处理的查询请求数目远低于OLTP系统，但每个查询通常要求非常苛刻，需要在短时间内扫描数百万条记录。磁盘带宽（不是寻道时间）通常是瓶颈，而面向列的存储对于这种工作负载成为日益流行的解决方案。

在OLTP方面，由两个主要流派的存储引擎：

- 日志结构流派，它只允许追加式更新文件和删除过时的文件，但不会修改已写入的文件。BitCask、SSTables、LSM-tree、LevelDB、Cassandra、HBase、Lucene等属于此类。



- 原地更新流派，将磁盘视为可以覆盖的一组固定大小的页。B-tree是这一哲学的最典型代表，它已用于所有主要的关系数据库，以及大量的非关系数据库。

日志结构的存储引擎是一个相对较新的方案。其关键思想是系统地将磁盘上随机访问写入转为顺序写入，由于硬盘驱动器和SSD的性能特性，可以实现更高的写入吞吐量。

此外，简要介绍了一些更复杂的索引结构，以及为全内存而优化的数据库。

然后，从存储引擎的内部间接地探索了典型数据仓库的总体架构。由此说明为什么分析工作负载与OLTP如此不同：当查询需要在大量行中顺序扫描时，索引的关联性就会显著降低。相反，最重要的是非常紧凑地编码数据，以尽量减少磁盘读取的数据量。我们讨论了列存储如何帮助实现这一目标。

作为应用开发人员，掌握更多有关存储引擎内部的知识，可以更好地了解哪种工具最适合你的具体应用。如果还需要进一步调整数据库的可调参数，这些理解还可以帮助开发者正确评估调高或调低参数所带来的影响。

尽管本章不能让你成为某个特定存储引擎的调优专家，但希望帮你获得足够的知识与见解，以充分理解所选择的数据库。

## 参考文献

- [1] Alfred V. Aho, John E. Hopcroft, and Jeffrey D. Ullman: *Data Structures and Algorithms*. Addison-Wesley, 1983. ISBN: 978-0-201-00023-8.
- [2] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein: *Introduction to Algorithms*, 3rd edition. MIT Press, 2009. ISBN: 978-0-262-53305-8.
- [3] Justin Sheehy and David Smith: “Bitcask: A Log-Structured Hash Table for Fast Key/Value Data,” Basho Technologies, April 2010.
- [4] Yinan Li, Bingsheng He, Robin Jun Yang, et al.: “Tree Indexing on Solid State Drives,” *Proceedings of the VLDB Endowment*, volume 3, number 1, pages 1195-1206, September 2010.
- [5] Goetz Graefe: “Modern B-Tree Techniques,” *Foundations and Trends in Databases*, volume 3, number 4, pages 203-402, August 2011. doi:10.1561/19000000028.





- [6] Jeffrey Dean and Sanjay Ghemawat: “LevelDB Implementation Notes,” *leveldb.googlecode.com*.
- [7] Dhruba Borthakur: “The History of RocksDB,” *rocksdb.blogspot.com*, November 24, 2013.
- [8] Matteo Bertozzi: “Apache HBase I/O - HFile,” *blog.cloudera.com*, June, 29 2012.
- [9] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, et al.: “Bigtable: A Distributed Storage System for Structured Data,” at *7th USENIX Symposium on Operating System Design and Implementation (OSDI)*, November 2006.
- [10] Patrick O’Neil, Edward Cheng, Dieter Gawlick, and Elizabeth O’Neil: “The Log-Structured Merge-Tree (LSM-Tree),” *Acta Informatica*, volume 33, number 4, pages 351-385, June 1996. doi:10.1007/s002360050048
- [11] Mendel Rosenblum and John K. Ousterhout: “The Design and Implementation of a Log-Structured File System,” *ACM Transactions on Computer Systems*, volume 10, number 1, pages 26-52, February 1992. doi:10.1145/146941.146943.
- [12] Adrien Grand: “What Is in a Lucene Index?,” at *Lucene/Solr Revolution*, November 14, 2013.
- [13] Deepak Kandept: “Hacking Lucene-The Index Format,” *hackerlabs.org*, October 1, 2011.
- [14] Michael McCandless: “Visualizing Lucene’s Segment Merges,” *blog.mikemccandless.com*, February 11, 2011.
- [15] Burton H. Bloom: “Space/Time Trade-offs in Hash Coding with Allowable Errors,” *Communications of the ACM*, volume 13, number 7, pages 422-426, July 1970. doi:10.1145/362686.362692.
- [16] “Operating Cassandra: Compaction,” Apache Cassandra Documentation v4.0, 2016.
- [17] Rudolf Bayer and Edward M. McCreight: “Organization and Maintenance of Large Ordered Indices,” Boeing Scientific Research Laboratories, Mathematical and Information Sciences Laboratory, report no. 20, July 1970.



- [18] Douglas Comer: “The Ubiquitous B-Tree,” *ACM Computing Surveys*, volume 11, number 2, pages 121-137, June 1979. doi:10.1145/356770.356776.
- [19] Emmanuel Goossaert: “Coding for SSDs,” *codecapsule.com*, February 12, 2014.
- [20] C. Mohan and Frank Levine: “ARIES/IM: An Efficient and High Concurrency Index Management Method Using Write-Ahead Logging,” at *ACM International Conference on Management of Data (SIGMOD)*, June 1992. doi: 10.1145/130283.130338.
- [21] Howard Chu: “LDAP at Lightning Speed,” at *Build Stuff '14*, November 2014.
- [22] Bradley C. Kuszmaul: “A Comparison of Fractal Trees to Log-Structured Merge (LSM) Trees,” *tokutek.com*, April 22, 2014.
- [23] Manos Athanassoulis, Michael S. Kester, Lukas M. Maas, et al.: “Designing Access Methods: The RUM Conjecture,” at *19th International Conference on Extending Database Technology (EDBT)*, March 2016. doi:10.5441/002/edbt.2016.42.
- [24] Peter Zaitsev: “InnoDB Double Write,” *percona.com*, August 4, 2006.
- [25] Tomas Vondra: “On the Impact of Full-Page Writes,” *blog.2ndquadrant.com*, November 23, 2016.
- [26] Mark Callaghan: “The Advantages of an LSM vs a B-Tree,” *smalldatum.blogspot.co.uk*, January 19, 2016.
- [27] Mark Callaghan: “Choosing Between Efficiency and Performance with RocksDB,” at *Code Mesh*, November 4, 2016.
- [28] Michi Mutsuzaki: “MySQL vs. LevelDB,” *github.com*, August 2011.
- [29] Benjamin Coverston, Jonathan Ellis, et al.: “CASSANDRA-1608: Redesigned Compaction, *issues.apache.org*, July 2011.
- [30] Igor Canadi, Siying Dong, and Mark Callaghan: “RocksDB Tuning Guide,” *github.com*, 2016.
- [31] MySQL 5.7 *Reference Manual*. Oracle, 2014.
- [32] *Books Online for SQL Server 2012*. Microsoft, 2012.
- [33] Joe Webb: “Using Covering Indexes to Improve Query Performance,” *simpletalk*.



com, 29 September 2008.

[34] Frank Ramsak, Volker Markl, Robert Fenk, et al.: “Integrating the UB-Tree into a Database System Kernel,” at *26th International Conference on Very Large Data Bases (VLDB)*, September 2000.

[35] The PostGIS Development Group: “PostGIS 2.1.2dev Manual,” *postgis.net*, 2014.

[36] Robert Escriva, Bernard Wong, and Emin Gün Sirer: “HyperDex: A Distributed, Searchable Key-Value Store,” at *ACM SIGCOMM Conference*, August 2012. doi: 10.1145/2377677.2377681

[37] Michael McCandless: “Lucene’s FuzzyQuery Is 100 Times Faster in 4.0,” *blog.mikemccandless.com*, March 24, 2011.

[38] Steffen Heinz, Justin Zobel, and Hugh E. Williams: “Burst Tries: A Fast, Efficient Data Structure for String Keys,” *ACM Transactions on Information Systems*, volume 20, number 2, pages 192–223, April 2002. doi:10.1145/506309.506312.

[39] Klaus U. Schulz and Stoyan Mihov: “Fast String Correction with Levenshtein Automata,” *International Journal on Document Analysis and Recognition*, volume 5, number 1, pages 67–85, November 2002. doi:10.1007/s10032-002-0082-8

[40] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze: *Introduction to Information Retrieval*. Cambridge University Press, 2008. ISBN: 978-0-521-86571-5, available online at [nlp.stanford.edu/IR-book](http://nlp.stanford.edu/IR-book).

[41] Michael Stonebraker, Samuel Madden, Daniel J. Abadi, et al.: “The End of an Architectural Era (It’s Time for a Complete Rewrite),” at *33rd International Conference on Very Large Data Bases (VLDB)*, September 2007.

[42] “VoltDB Technical Overview White Paper,” VoltDB, 2014.

[43] Stephen M. Rumble, Ankita Kejriwal, and John K. Ousterhout: “Log-Structured Memory for DRAM-Based Storage,” at *12th USENIX Conference on File and Storage Technologies (FAST)*, February 2014.

[44] Stavros Harizopoulos, Daniel J. Abadi, Samuel Madden, and Michael Stonebraker: “OLTP Through the Looking Glass, and What We Found There,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2008. doi:





10.1145/1376616.1376713.

[45] Justin DeBrabant, Andrew Pavlo, Stephen Tu, et al.: “Anti-Caching: A New Approach to Database Management System Architecture,” *Proceedings of the VLDB Endowment*, volume 6, number 14, pages 1942–1953, September 2013.

[46] Joy Arulraj, Andrew Pavlo, and Subramanya R. Dulloor: “Let’s Talk About Storage & Recovery Methods for Non-Volatile Memory Database Systems,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2015. doi: 10.1145/2723372.2749441.

[47] Edgar F. Codd, S. B. Codd, and C. T. Salley: “Providing OLAP to User-Analysts: An IT Mandate,” E. F. Codd Associates, 1993.

[48] Surajit Chaudhuri and Umeshwar Dayal: “An Overview of Data Warehousing and OLAP Technology,” *ACM SIGMOD Record*, volume 26, number 1, pages 65–74, March 1997. doi:10.1145/248603.248616.

[49] Per-Åe Larson, Cipri Clinciu, Campbell Fraser, et al.: “Enhancements to SQL Server Column Stores,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2013.

[50] Franz Fäber, Norman May, Wolfgang Lehner, et al.: “The SAP HANA Database- An Architecture Overview,” *IEEE Data Engineering Bulletin*, volume 35, number 1, pages 28–33, March 2012.

[51] Michael Stonebraker: “The Traditional RDBMS Wisdom Is (Almost Certainly) All Wrong,” presentation at EPFL, May 2013.

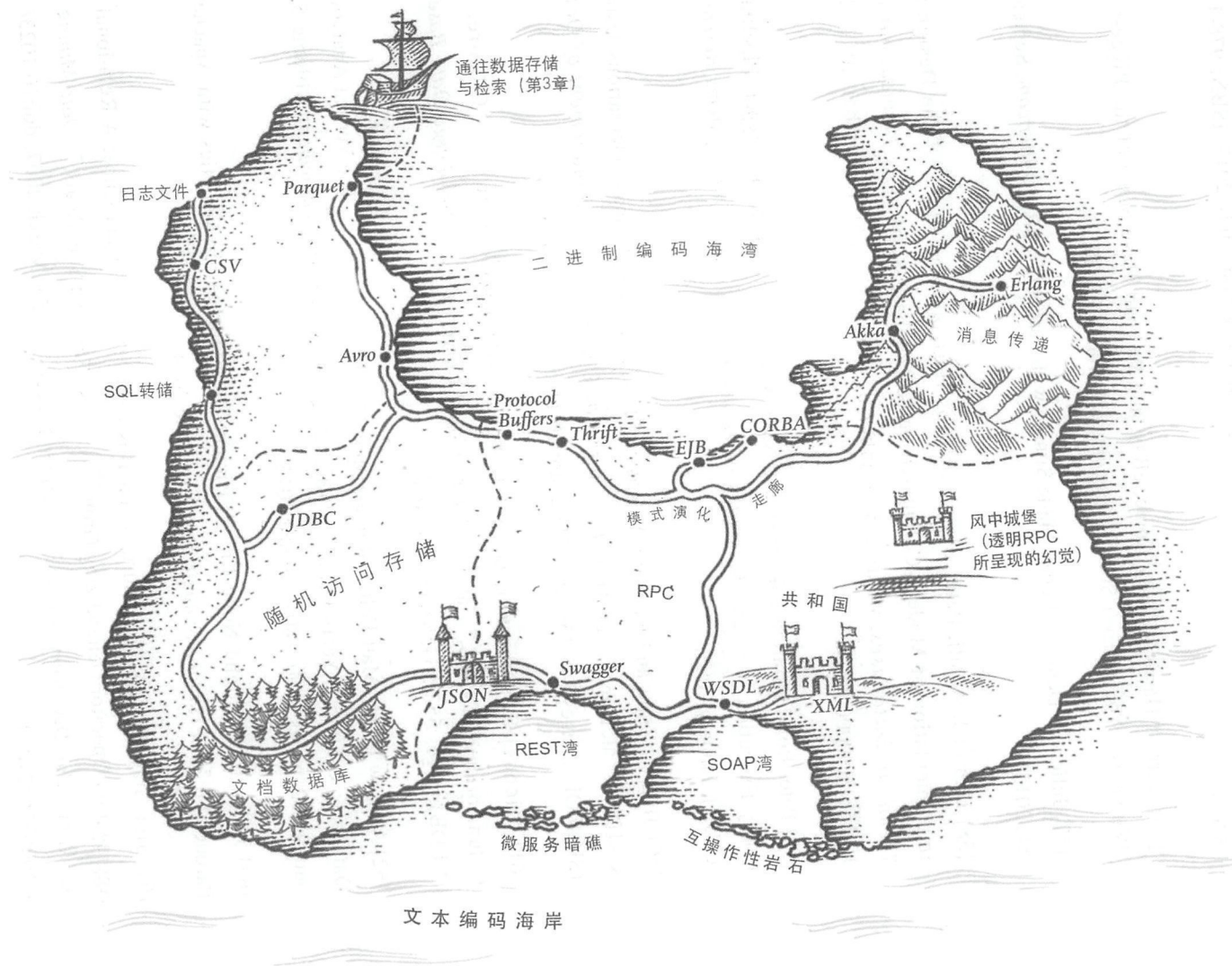
[52] Daniel J. Abadi: “Classifying the SQL-on-Hadoop Solutions,” *hadapt.com*, October 2, 2013.

[53] Marcel Kornacker, Alexander Behm, Victor Bittorf, et al.: “Impala: A Modern, Open-Source SQL Engine for Hadoop,” at *7th Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2015.

[54] Sergey Melnik, Andrey Gubarev, Jing Jing Long, et al.: “Dremel: Interactive Analysis of Web-Scale Datasets,” at *36th International Conference on Very Large Data Bases (VLDB)*, pages 330–339, September 2010.



- [55] Ralph Kimball and Margy Ross: *The Data Warehouse Toolkit: The Definitive Guide to Dimensional Modeling*, 3rd edition. John Wiley & Sons, July 2013. ISBN: 978-1-118-53080-1.
- [56] Derrick Harris: “Why Apple, eBay, and Walmart Have Some of the Biggest Data Warehouses You’ve Ever Seen,” *gigaom.com*, March 27, 2013.
- [57] Julien Le Dem: “Dremel Made Simple with Parquet,” *blog.twitter.com*, September 11, 2013.
- [58] Daniel J. Abadi, Peter Boncz, Stavros Harizopoulos, et al.: “The Design and Implementation of Modern Column-Oriented Database Systems,” *Foundations and Trends in Databases*, volume 5, number 3, pages 197–280, December 2013. doi: 10.1561/19000000024.
- [59] Peter Boncz, Marcin Zukowski, and Niels Nes: “MonetDB/X100: Hyper- Pipelining Query Execution,” at *2nd Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2005.
- [60] Jingren Zhou and Kenneth A. Ross: “Implementing Database Operations Using SIMD Instructions,” at *ACM International Conference on Management of Data (SIGMOD)*, pages 145-156, June 2002. doi:10.1145/564691.564709.
- [61] Michael Stonebraker, Daniel J. Abadi, Adam Batkin, et al.: “C-Store: A Columnoriented DBMS,” at *31st International Conference on Very Large Data Bases (VLDB)*, pages 553-564, September 2005.
- [62] Andrew Lamb, Matt Fuller, Ramakrishna Varadarajan, et al.: “The Vertica Analytic Database: C-Store 7 Years Later,” *Proceedings of the VLDB Endowment*, volume 5, number 12, pages 1790–1801, August 2012.
- [63] Julien Le Dem and Nong Li: “Efficient Data Storage for Analytics with Apache Parquet 2.0,” at *Hadoop Summit*, San Jose, June 2014.
- [64] Jim Gray, Surajit Chaudhuri, Adam Bosworth, et al.: “Data Cube: A Relational Aggregation Operator Generalizing Group-By, Cross-Tab, and Sub-Totals,” *Data Mining and Knowledge Discovery*, volume 1, number 1, pages 29–53, March 2007. doi:10.1023/A:1009726021843.







# 数据编码与演化

一切都在改变，一刻都没有停止。

——Heraclitus，如柏拉图在Cratylus引用的那样（公元前360年）

应用程序不可避免地需要随时间而变化、调整。当新产品推出，或为了更好地理解用户需求，或商业环境发生变化时，就需要不断地添加或修改功能。在第1章中我们介绍了可演化的思想，提出应该构建可适应变化的系统（参阅第1章“可演化性：易于改变”）。

在大多数情况下，更改应用程序功能时，也需要更改其存储的数据：可能需要捕获新的字段或记录类型，或者需要以新的方式呈现已有数据。

第2章所讨论的数据模型有不同的方法来应对这种变化。关系数据库通常假设数据库中的所有数据都符合一种模式，尽管该模式可以改变（通过模式迁移，即ALTER语句），这样在任何一个给定时间点都只有一个有效的模式。相比之下，读时模式（“无模式”）数据库不强制执行模式，所以数据库包含了不同时间写入的新旧数据的混合体（参阅第2章“文档模型中的模式灵活性”）。

当数据格式或模式发生变化时，经常需要对应用程序代码进行相应的调整（例如，向记录中添加新字段，然后应用程序代码开始读取和写入该字段）。然而，对于一个大型应用系统，代码更迭往往并非易事：

- 对于服务器端应用程序，可能需要执行滚动升级（也被称为分阶段发布），每次将新版本部署到少数几个节点，检查新版本是否正常运行，然后逐步在所有节点上升级新的代码。这样新版本部署无需服务暂停，从而支持更频繁的版本发布和



更好的演化。

- 对于客户端应用程序，只能寄望于用户，然而他们在一段时间内可能不会马上安装更新。

这意味着新旧版本的代码，以及新旧数据格式，可能会同时在系统内共存。为了使系统继续顺利运行，需要保持双向的兼容性：

#### 向后兼容

较新的代码可以读取由旧代码编写的的数据。

#### 向前兼容

较旧的代码可以读取由新代码编写的的数据。

向后兼容通常不难实现：作为新代码的作者，清楚旧代码所编写的的数据格式，因此可以比较明确地处理这些旧数据（如果需要，只需保留旧的代码来读取旧的数据）。向前兼容可能会比较棘手，它需要旧代码忽略新版本的代码所做的添加。

在本章中，将介绍多种编码数据的格式，包括JSON、XML、Protocol Buffers、Thrift和Avro。特别地，我们将讨论它们如何处理模式变化，以及如何支持新旧数据和新旧代码共存的系统。之后，还将讨论这些格式如何用于数据存储和通信场景，包括在Web服务中，具象状态传输（Representational State Transfer, REST）和远程过程调用（remote procedure calls, RPC）以及消息传递系统，如actors和消息队列。

## 数据编码格式

程序通常使用（至少）两种不同的数据表示形式：

1. 在内存中，数据保存在对象、结构体、列表、数组、哈希表和树等结构中。这些数据结构针对CPU的高效访问和操作进行了优化（通常使用指针）。
2. 将数据写入文件或通过网络发送时，必须将其编码为某种自包含的字节序列（例如JSON文档）。由于指针对其他进程没有意义，所以这个字节序列表示看起来与内存中使用的数据结构大不一样<sup>注1</sup>。

因此，在这两种表示之间需要进行类型的转化。从内存中的表示到字节序列的转化称为编码（或序列化等），相反的过程称为解码（或解析，反序列化）<sup>注2</sup>。

注1：除了某些特殊情况，例如某些内存映射文件或直接在压缩数据上操作（如第3章的“列压缩”部分所述）。

注2：请注意，编码与加密无关。加密不在讨论范围之内。



### 术语冲突

不幸的是，序列化一词也用于事务处理（参看第7章），且意义完全不同。为了避免在本书中产生冲突，尽管序列化可能是更常见的术语，我们仍坚持使用编码。

由于这是一个常见的问题，因此存在许多不同的库和编码格式可供选择。我们在此做一个简要概述。

## 语言特定的格式

许多编程语言都内置支持将内存中的对象编码为字节序列。例如，Java有`java.io.Serializable`<sup>[1]</sup>，Ruby有`Marshal`<sup>[2]</sup>，Python有`pickle`<sup>[3]</sup>等。此外，还有许多第三方库，例如用于Java的`Kryo`<sup>[4]</sup>。

这些编码库使用起来非常方便，它们只需要很少的额外代码即可保存或恢复内存中的对象。然而，这里也有一些深层次的问题：

- 编码通常与特定的编程语言绑定在一起，而用另一种语言访问数据就非常困难。如果用这种编码方式存储或传输数据，可能在很长时间内须使用当前的编程语言，并且不能将系统与其他组织（可能使用不同的语言）的系统方便地集成在一起。
- 为了在相同的对象类型中恢复数据，解码过程需要能够实例化任意的类。这经常导致一些安全问题<sup>[5]</sup>：如果攻击者可以让应用程序解码任意的字节序列，那么它们可以实例化任意的类，这通常意味着，它们可以做一些可怕的事情，比如远程执行任意代码<sup>[6, 7]</sup>。
- 在这些库中，多版本数据通常是次要的，主要目标是快速且简单地编码数据，所以它们经常忽略向前和向后兼容性问题。
- 效率（编码或解码花费的CPU时间，以及编码结构的大小）通常也是次要的。例如，Java的内置序列化由于其糟糕的性能和臃肿的编码而广为诟病<sup>[8]</sup>。

由于这些原因，使用语言内置的编码方案通常不是个好主意，除非只是为了临时尝试。

## JSON、XML与二进制变体

目标转向可由不同编程语言编写和读取的标准化编码，显然JSON和XML是其中佼佼者。它们广为人知，得到广泛的支持，以及同样广泛的批评。XML经常被批评过于





冗长和不必要的复杂<sup>[9]</sup>。JSON受欢迎主要是由于它在Web浏览器中内置支持（因为是JavaScript的一个子集）以及相对于XML的简单性。CSV是另一种流行的与语言无关的格式，尽管功能较弱。

JSON、XML和CSV都是文本格式，因此具有不错的可读性（尽管语法容易引发争论）。除了表面的语法问题之外，它们也有一些微妙的问题：

- 数字编码有很多模糊之处。在XML和CSV中，无法区分数字和碰巧由数字组成的字符串（除了引用外部模式）。JSON区分字符串和数字，但不区分整数和浮点数，并且不指定精度。

这在处理大数字时是一个问题，大于 $2^{53}$ 的整数在IEEE 754双精度浮点数中不能精确表示，所以这些数字在使用浮点数（如JavaScript）的语言中进行分析时，会变得不准确。Twitter上有一个大于 $2^{53}$ 的数字的例子，它使用一个64位的数字来标识每条推文。Twitter的API返回的JSON包含两次推特ID，一次是JSON数字，一次是十进制字符串，以解决JavaScript应用程序没有正确解析数字的问题<sup>[10]</sup>。

- JSON和XML对Unicode字符串（即人类可读文本）有很好的支持，但是它们不支持二进制字符串（没有字符编码的字节序列）。二进制字符串是一个有用的功能，所以人们通过使用Base64将二进制数据编码为文本来解决这个限制。然后，模式可以表明该值应该被解释为Base64编码。虽然可行，但有点混乱，并且数据大小增加了33%。
- XML<sup>[11]</sup>和JSON<sup>[12]</sup>都有可选的模式支持。这些模式语言相当强大，因此学习和实现起来也比较复杂。XML模式的使用相当广泛，但许多基于JSON的工具并不局限于使用模式。由于数据（例如数字和二进制字符串）的正确解释取决于模式中的信息，因此不使用XML/JSON架构的应用程序可能不得不硬编码适当的编码/解码逻辑。
- CSV没有任何模式，因此应用程序需要定义每行和每列的含义。如果应用程序更改添加新的行或列，则必须手动处理该更改。CSV也是一个相当模糊的格式（如果一个值包含逗号或换行符，会发生什么）。尽管其转义规则已经被正式指定<sup>[13]</sup>，但并不是所有的解析器都能正确地实现它们。

尽管存在这些或那些缺陷，但JSON、XML和CSV已经可用于很多应用。特别是作为数据交换格式（即将数据从一个组织发送到另一个组织），它们非常受欢迎。在这些情况下，只要人们就格式本身达成一致，格式多么美观或者高效往往不太重要。让不同的组织达成格式一致的难度通常超过了所有其他问题。



## 二进制编码

对于仅在组织内部使用的数据，使用最小公分母编码格式则较为顺畅。例如，可以选择更紧凑或更快的解析格式。对于一个小数据集来说，收益可以忽略不计，但一旦达到TB级别，数据格式的选择就会产生很大的影响。

JSON不像XML那么冗长，但与二进制格式相比，两者仍然占用大量空间。这种观察导致开发了大量的二进制编码，用以支持JSON（举几个例子，如MessagePack、BSON、BJSON、UBJSON、BISON和Smile）和XML（如WBXML和Fast Infoset）。这些格式已经被很多细分领域所采用，但是没有一个像JSON和XML那样被广泛采用。

其中一些格式还扩展了数据类型集（例如，区分整数和浮点数，或者增加对二进制字符串的支持），但其他格式保持JSON / XML数据模型不变。特别是，由于它们没有规定模式，所以需要在编码数据时包含所有的对象字段名称。也就是说，在示例4-1中的JSON文档的二进制编码中，它们将需要在某处包含字符串userName、favoriteNumber和interest。

示例4-1：一条样本记录，本章后续的多种二进制格式编码都以该记录来示例说明

```
{
  "userName": "Martin",
  "favoriteNumber": 1337,
  "interests": ["daydreaming", "hacking"]
}
```

来看一个MessagePack的例子，它是一种JSON的二进制编码。图4-1展示了采用MessagePack<sup>[14]</sup>对示例4-1的JSON文档进行编码所得到的字节序列。前几个字节如下：

1. 第一个字节0x83，表示接下来是包含三个字段（最低四位=0x03）的对象（最高四位=0x80）（如果想知道当对象的字段数超过15个，4bit已经无法容纳，会发生什么情况。结果是，它会得到一个不同的类型指示符，并且字段数被编码为两个或四个字节）。
2. 第二个字节0xa8，表示接下来是八字节长的字符串（最高四位=0xa0，最低四位=0x08）。
3. 再往下的八字节是ASCII中的字段名称userName。由于之前已经指出了长度，因此不需要任何标记来告诉字符串结束的位置（或任何转义）。
4. 接下来的七字节使用前缀0xa6对6个字母的字符串值Martin进行编码，依此类推。

二进制编码的长度为66字节，仅略小于文本JSON编码（去掉空格）占用的81字节。JSON的所有二进制编码在这方面是相似的。目前对于如此小的空间缩减（也许解析



速度可以加快) 是否值得失去可读性仍存有疑问。

在下面的章节中, 我们将看到如何做得更好, 只用32字节即可完成对同样的记录进行编码。

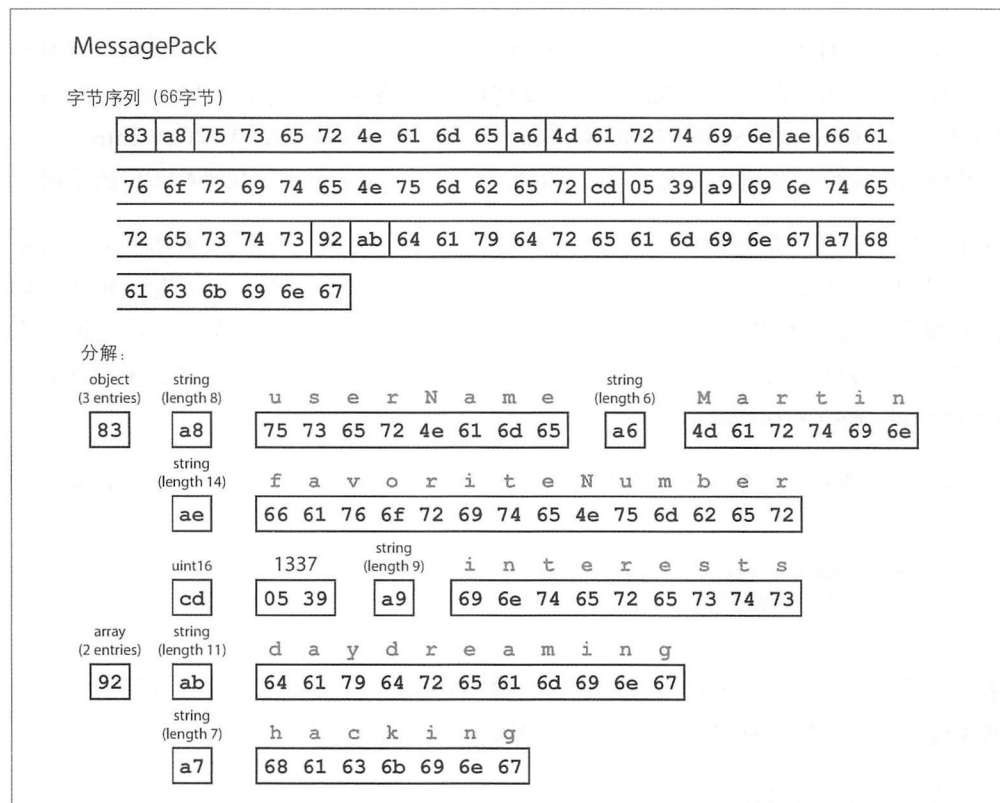


图4-1: 用MessagePack来编码的样本记录 (示例4-1)

## Thrift与Protocol Buffers

Apache Thrift<sup>[15]</sup>和Protocol Buffers (protobuf)<sup>[16]</sup>是基于相同原理的两种二进制编码库。Protocol Buffers最初是在Google开发的, Thrift最初是在Facebook开发的, 并且都是在2007~2008年开源的<sup>[17]</sup>。

Thrift和Protocol Buffers都需要模式来编码任意的数据。为了用Thrift对示例4-1中的数据进行编码, 可以使用Thrift接口定义语言 (IDL) 来描述模式, 如下所示:

```
struct Person {
    1: required string    userName,
    2: optional i64      favoriteNumber,
```





```
3: optional list<string> interests
}
```

Protocol Buffers的等价模式定义看起来非常相似：

```
message Person {
  required string user_name      = 1;
  optional int64 favorite_number = 2;
  repeated string interests      = 3;
}
```

Thrift和Protocol Buffers各有对应的代码生成工具，采用和上面类似的模式定义，并生成支持多种编程语言<sup>[18]</sup>的类。应用程序可以直接调用生成的代码来编码或解码该模式的记录。

用这个模式编码的数据是什么样的？令人困惑的是，Thrift有两种不同的二进制编码格式<sup>注3</sup>，分别称为BinaryProtocol和CompactProtocol。先来看看BinaryProtocol，以这种格式编码示例4-1需要59字节，如图4-2所示<sup>[19]</sup>。

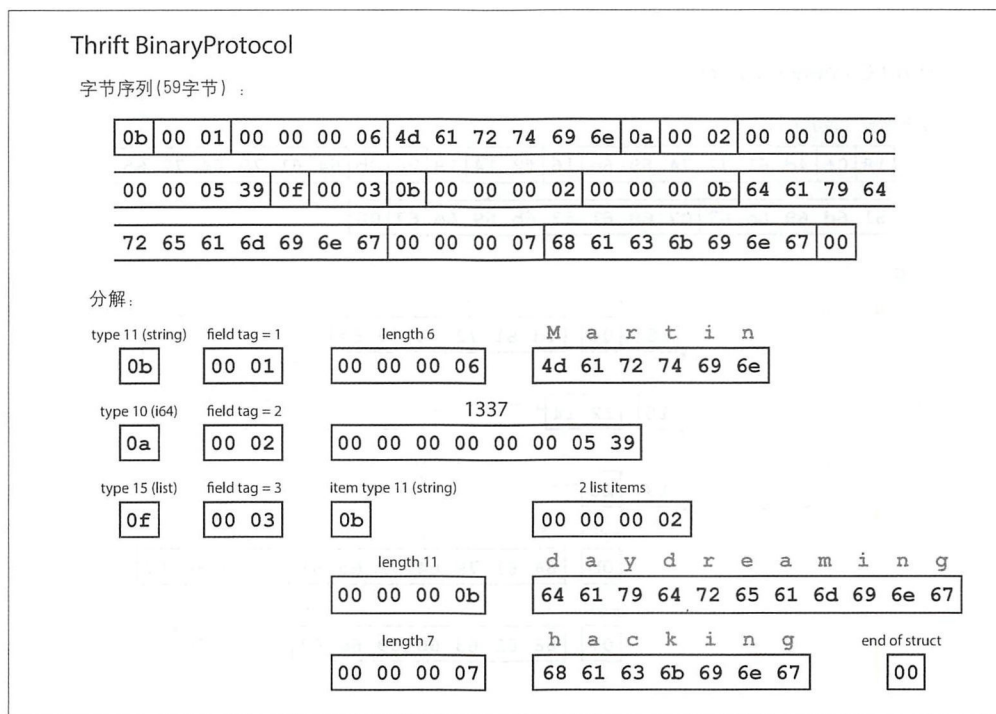


图4-2：使用Thrift的BinaryProtocol编码的示例记录

注3： 实际上，它具有三种：BinaryProtocol、CompactProtocol和DenseProtocol。虽然DenseProtocol只支持C++实现，也正因如此，所以它不算跨语言<sup>[18]</sup>。除此之外，它还有两种不同的基于JSON的编码格式<sup>[19]</sup>。太有趣了！



与图4-1类似，每个字段都有一个类型注释（用于指示它是否是字符串、整数、列表等），并且可以在需要时指定长度（包括字符串的长度、列表中的项数）。与之前类似，数据中出现的字符串（“Martin”，“daydreaming”，“hacking”）也被编码为ASCII（或者更确切地说，UTF-8）。

与图4-1相比，最大的区别是没有字段名（userName、favoriteNumber和interest）。相反，编码数据包含数字类型的字段标签（1、2和3）。这些是模式定义中出现的数字。字段标签就像字段的别名，用来指示当前的字段，但更为紧凑，可以省去引用字段全名。

Thrift CompactProtocol编码在语义上等同于BinaryProtocol，但如图4-3所示，它将相同的信息打包成只有34字节。它通过将字段类型和标签号打包到单字节中，并使用可变长度整数来实现。对数字1337，不使用全部8字节，而是使用两个字节进行编码，每字节的最高位用来指示是否还有更多的字节。这意味着-64~63之间的数字被编码为一字节，-8192~8191之间的数字被编码成两个字节等。更大的数字需要更多字节。

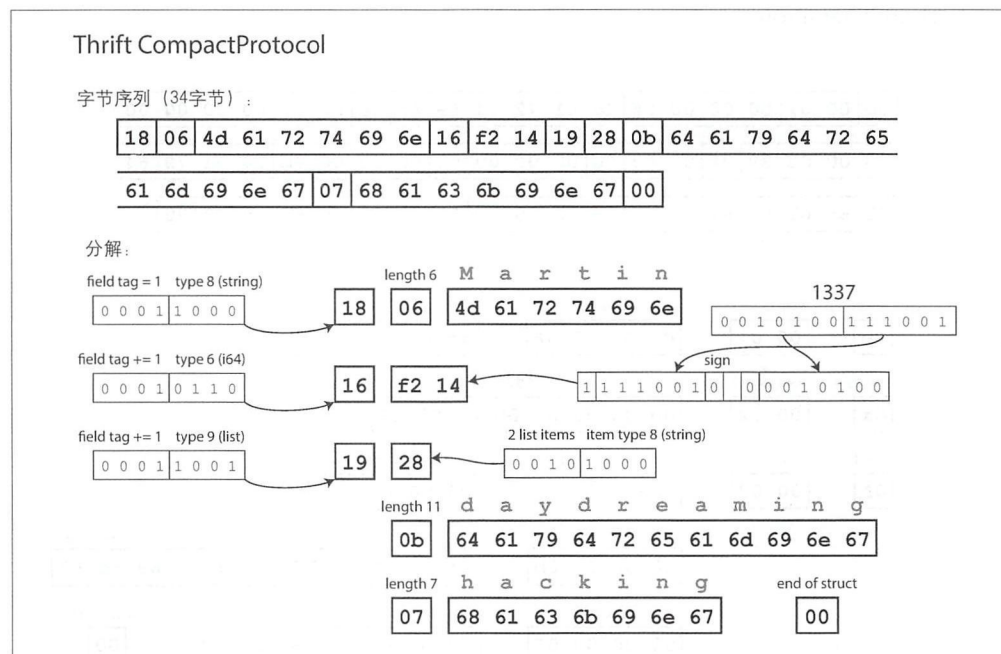


图4-3：采用Thrift CompactProtocol来编码的样本记录

最后，Protocol Buffers（只有一种二进制编码格式）对相同的数据进行编码，如图4-4所示。它的位打包方式略有不同，但与Thrift的CompactProtocol非常相似。Protocol Buffers只用33字节可以表示相同的记录。

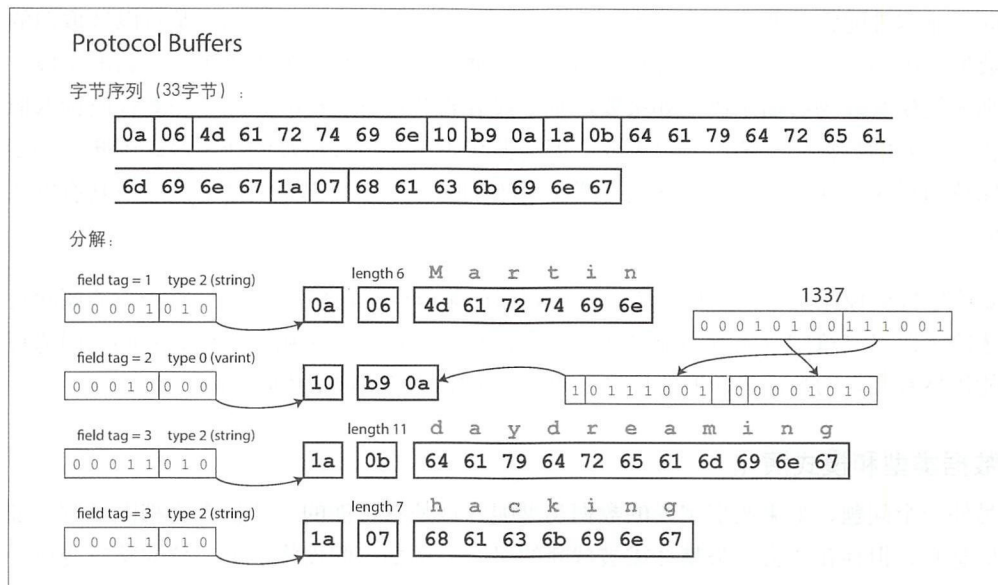


图4-4: 使用Protocol Buffers编码的示例记录

需要注意的一个细节是，在前面所示的模式中，每个字段被标记为required（必须）或optional（可选），但这对字段如何编码没有任何影响（二进制数据中不会指示某字段是否必须）。区别在于，如果字段设置了required，但字段未填充，则运行时检查将出现失败，这对于捕获错误非常有用。

## 字段标签和模式演化

之前说过，模式不可避免地需要随着时间而不断变化，称之为模式演化。那么Thrift和Protocol Buffers如何在保持向后和向前兼容性的同时应对模式更改呢？

从示例中可以看到，一条编码记录只是一组编码字段的拼接。每个字段由其标签号（示例模式中的数字1、2、3）标识，并使用数据类型（例如字符串或整数）进行注释。如果没有设置字段值，则将其从编码的记录中简单地忽略。由此可以看出，字段标签（field tag）对编码数据的含义至关重要。可以轻松更改模式中字段的名称，而编码永远不直接引用字段名称。但不能随便更改字段的标签，它会导致所有现有编码数据无效。

可以添加新的字段到模式，只要给每个字段一个新的标记号码。如果旧的代码（不知道添加的新标记号码）试图读取新代码写入的数据，包括一个它不能识别的标记号码中新的字段，则它可以简单地忽略该字段。实现时，通过数据类型的注释来通知解析器跳过特定的字节数。这样可以实现向前兼容性，即旧代码可以读取由新代码编写的记录。





向后兼容性呢？只要每个字段都有一个唯一的标记号码，新的代码总是可以读取旧的数据，因为标记号码仍然具有相同的含义。唯一的细节是，如果添加一个新的字段，则无法使其成为必需字段。如果要添加字段并将其设置为required，当新代码读取旧代码写入的数据，则该检查将失败，因为旧代码不会写入添加的新字段。因此，为了保持向后兼容性，在模式的初始部署之后添加的每个字段都必须都是可选的或具有默认值。

删除字段就像添加字段一样，不过向后和向前兼容性问题相反。这意味着只能删除可选的字段（必填字段永远不能被删除），而且不能再次使用相同的标签号码（因为可能仍然有写入的数据包含旧的标签号码，而该字段必须被新代码忽略）。

## 数据类型和模式演化

另外一个问题，如果改变字段的数据类型呢？这是有可能的（请检查文档以了解详细信息），但存在值会丢失精度或被截断的风险。例如，假设将一个32位的整数变成一个64位的整数。新代码可以容易地读取旧代码写入的数据，因为解析器可以用零填充任何缺失的位。但是，如果旧代码读取新代码写入的数据，旧代码仍然使用32位变量来保存该值。如果解码的64位值不适合32位，则它将被截断。

Protocol Buffers的一个奇怪的细节是，它没有列表或数组数据类型，而是有字段的重复标记（repeated，这是必需和可选之外的第三个选项）。如图4-4所示，对于重复字段，表示同一个字段标签只是简单地多次出现在记录中。可以将可选（单值）字段更改为重复（多值）字段。读取旧数据的新代码会看到一个包含零个或一个元素的列表（取决于该字段是否存在）。读取新数据的旧代码只能看到列表的最后一个元素。

Thrift有专用的列表数据类型，它使用列表元素的数据类型进行参数化。它不支持Protocol Buffers那样从单值到多值的改变，但是它具有支持嵌套列表的优点。

## Avro

Apache Avro<sup>[20]</sup>是另一种二进制编码格式，它与Protocol Buffers和Thrift有着一些有趣的差异。由于Thrift不适合Hadoop的用例<sup>[21]</sup>，因此Avro在2009年作为Hadoop的子项目而启动。

Avro也使用模式来指定编码的数据结构。它有两种模式语言：一种（Avro IDL）用于人工编辑，另一种（基于JSON）更易于机器读取。

用Avro IDL编写的示例模式如下所示：



```
record Person {
  string      userName;
  union { null, long } favoriteNumber = null;
  array<string> interests;
}
```

该模式的等价JSON表示如下:

```
{
  "type": "record",
  "name": "Person",
  "fields": [
    { "name": "userName", "type": "string" },
    { "name": "favoriteNumber", "type": ["null", "long"], "default": null },
    { "name": "interests", "type": { "type": "array", "items": "string" } }
  ]
}
```

首先, 请注意模式中没有标签编号。如果使用这个模式编码示例记录(示例4-1), Avro二进制编码只有32字节长, 这是所见到的所有编码中最紧凑的。编码字节序列的分解如图4-5所示。

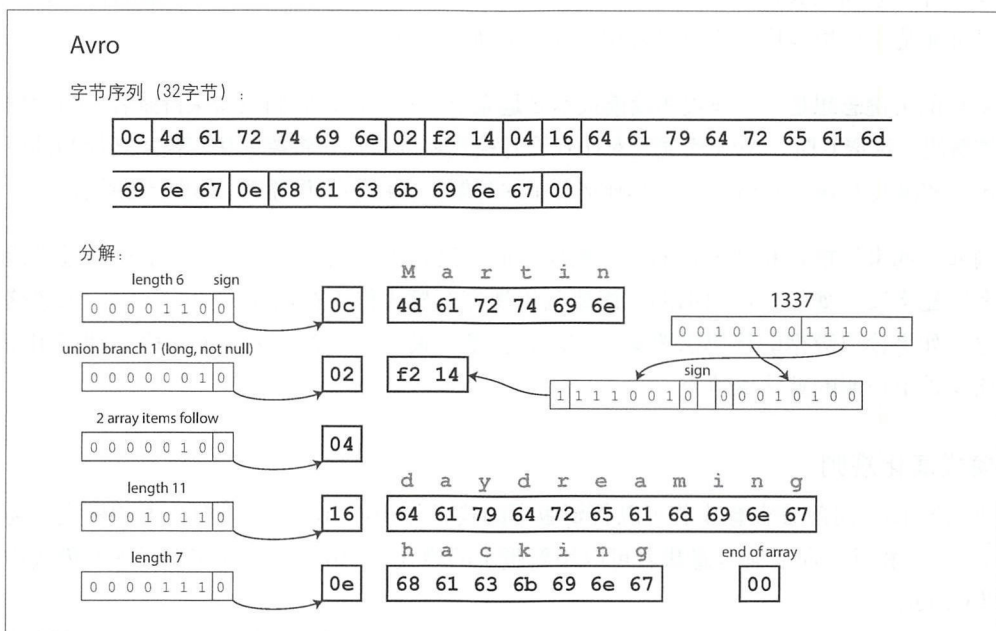


图4-5: 采用Avro编码的示例记录

如果检查图4-5的字节序列, 可以看到没有什么可以标识字段或数据类型。编码只是由连在一起的一些列值组成。一个字符串只是一个长度前缀, 后跟UTF-8字节流, 但



编码数据中没有任何内容告诉你它是一个字符串。它也可以是一个整数，或者其他什么类型。整数使用可变长度编码（与Thrift的CompactProtocol相同）进行编码。

为了解析二进制数据，按照它们出现在模式中的顺序遍历这些字段，然后直接采用模式告诉你每个字段的数据类型。这意味着，只有当读取数据的代码使用与写入数据的代码完全相同的模式时，才能正确解码二进制数据。读和写的模式如果有任何不匹配都将无法解码数据。

那么，Avro如何支持模式演化？

## 写模式与读模式

有了Avro，当应用程序想要对某些数据进行编码（例如将其写入文件或数据库，以及通过网络发送）时，它使用所知道的模式的任何版本来编码数据，例如，可以编译到应用程序中的模式。这被称为写模式。

当应用程序想要解码某些数据（例如从文件或数据库读取数据，或者从网络接收数据等）时，它期望数据符合某个模式，即读模式。这是应用程序代码所依赖的模式，代码可能是在应用程序的构建过程中基于模式而动态生成。

Avro的关键思想是，写模式和读模式不必是完全一模一样，它们只需保持兼容。当数据被解码（读取）时，Avro库通过对比查看写模式和读模式并将数据从写模式转换为读模式来解决其差异。Avro规范<sup>[20]</sup>明确定义了这种解决方法的工作原理，如图4-6所示。

例如，如果写模式和读模式的字段顺序不同，这也没有问题，因为模式解析通过字段名匹配字段。如果读取数据的代码遇到出现在写模式但不在读模式中的字段，则忽略它。如果读取数据的代码需要某个字段，但是写模式不包含该名称的字段，则使用在读模式中声明的默认值填充。

## 模式演化规则

使用Avro，向前兼容意味着可以将新版本的模式作为writer，并将旧版本的模式作为reader。相反，向后兼容意味着可以用新版本的模式作为reader，并用旧版本的模式作为writer。

为了保持兼容性，只能添加或删除具有默认值的字段（在Avro模式中，字段favourNumber的默认值为null）。例如，假设添加了一个带有默认值的字段，则此新字段存在于新模式中，而不是旧模式中。当使用新模式的reader读取使用旧模式写入的记录时，将为缺少的字段填充默认值。



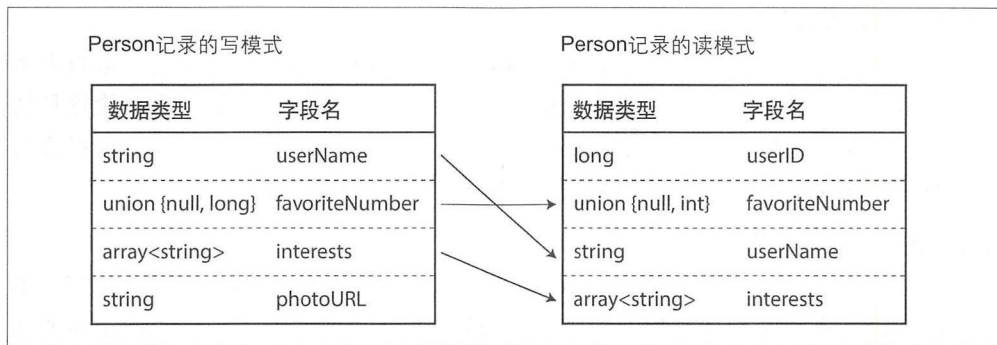


图4-6：由Avro读端解决写模式和读模式之间的差异

如果要添加一个没有默认值的字段，新的reader将无法读取旧的writer写的的数据，因此将破坏向后兼容性。如果要删除没有默认值的字段，旧reader将无法读取新writer写入的数据，因此将破坏向前兼容性。

在某些编程语言中，null是所有变量可以接受的默认值。但在Avro中并非如此：如果要允许字段为null，则必须使用联合类型。例如，`union{ null, long, string}`字段；表示该字段可以是数字、字符串或null。只有当null是联合的分支之一时，才可以使用它作为默认值<sup>注4</sup>。这比默认情况下所有类型都可为空显得更加冗长一些，但是通过明确什么能为null和不能为null可以帮助防止一些错误<sup>[22]</sup>。

因此，Avro不像Protocol Buffers和Thrift那样具有可选和必需的标签（而是有联合类型和默认值）。

只要Avro可以转换类型，就可以改变字段的数据类型。更改字段的名称也是可能的，但有点棘手：reader的模式可以包含字段名称的别名，因此它可以将旧writer模式字段名称与别名进行匹配。这意味着更改字段名称是向后兼容的，但不能向前兼容。同样，向联合类型添加分支也是向后兼容的，但不能向前兼容。

### 那么writer模式又是什么？

到目前为止，忽略了一个重要的问题：reader如何知道特定的数据采用哪个writer的模式编码的？在每个记录中都包含整个模式不太现实，因为模式有时甚至比编码数据还要大得多，这样二进制编码所节省的空间都变得没有意义。

答案取决于Avro使用的上下文。举几个例子：

注4：确切地说，默认值必须是联合的第一个可能类型，尽管这是Avro的特定限制，而不是联合类型的一般特征。



## 有很多记录的大文件

Avro的一个常见用途，尤其是在Hadoop的上下文中，是用于存储包含数百万条记录的大文件，所有记录都使用相同的模式进行编码（将在第10章讨论这种情况）。在这种情况下，该文件的writer可以仅在文件的开头包含writer的模式信息。Avro通过指定一个文件格式（对象容器文件）来做到这一点。

## 具有单独写入记录的数据库

在数据库中，不同的记录可能在不同的时间点、使用不同的writer模式编写，不能假设所有记录都具有相同的模式。最简单的解决方案是在每个编码记录的开始处包含一个版本号，并在数据库中保留一个模式版本列表。reader可以获取记录，提取版本号，然后从数据库中查询该版本号的writer模式。使用该writer模式，它可以解码记录的其余部分。例如Espresso<sup>[23]</sup>就是这样工作的。

## 通过网络连接发送记录

当两个进程通过双向网络连接进行通信时，他们可以在建立连接时协商模式版本，然后在连接的生命周期中使用该模式。这也是Avro RPC协议（参阅本章后续的“基于服务的数据流：REST和RPC”）的基本原理。

在任何情况下，提供一个模式版本信息的数据库都非常有用，它可以充当一个说明文档来检查模式兼容性情况<sup>[24]</sup>。至于版本号，可以使用简单的递增整数，也可以使用对模式的哈希。

## 动态生成的模式

与Protocol Buffers和Thrift相比，Avro方法的一个优点是不包含任何标签号。为什么这很重要？在模式中保留一些数字有什么问题？

关键之处在于Avro对动态生成的模式更友好。例如，假如有一个关系数据库，想要把它的内容转储到一个文件中，并且希望使用二进制格式来避免上述文本格式（JSON、CSV、SQL）的问题。如果使用Avro，可以很容易地根据关系模式生成Avro模式，并使用该模式对数据库内容进行编码，然后将其全部转储到Avro对象容器文件<sup>[25]</sup>中。可以为每个数据库的表生成对应的记录模式，而每个列成为该记录中的一个字段。数据库中的列名称映射到Avro中的字段名称。

现在，如果数据库模式发生变化（例如，表中添加了一列，删除了一列），则可以从更新的数据库模式生成新的Avro模式，并用新的Avro模式导出数据。数据导出过程不需要关注模式的改变，每次运行时都可以简单地进行模式转换。任何读取新数据文件的人都会看到记录的字段已经改变，但是由于字段是通过名字来标识的，所以更新的



writer模式仍然可以与旧的reader模式匹配。

相比之下，如果使用Thrift或Protocol Buffers，则可能必须手动分配字段标签：每次数据库模式更改时，管理员都必须手动更新从数据库列名到字段标签的映射（这可能会自动化，但模式生成器必须非常小心，不要分配以前使用的字段标签）。这种动态生成的模式根本不是Thrift或Protocol Buffers的设计目标，而是Avro的设计目标。

## 代码生成和动态类型语言

Thrift和Protocol Buffers依赖于代码生成：在定义了模式之后，可以使用选择的编程语言生成实现此模式的代码。这在Java、C++或C#等静态类型语言中很有用，因为它允许使用高效的内存结构来解码数据，并且在编写访问数据结构的程序时，支持在IDE中进行类型检查和自动完成。

在动态类型编程语言中，如JavaScript、Ruby或Python，因为没有编译时类型检查，生成代码没有太多意义。代码生成在这些语言中经常被忽视，因为他们避免了明确的编译步骤。此外，对于动态生成的模式（例如从数据库表生成的Avro模式）的情况，代码生成对获取数据反而是不必要的障碍。

Avro为静态类型编程语言提供了可选的代码生成，但是它也可以在不生成代码的情况下直接使用。如果有一个对象容器文件（它嵌入了writer模式），可以简单地使用Avro库打开它，并用和查看JSON文件一样的方式查看数据。该文件是自描述的，它包含了所有必要的元数据。

此属性与动态类型数据处理语言（如Apache Pig<sup>[26]</sup>）结合使用时特别有用。在Apache Pig中，只需打开一些Avro文件，分析其内容，并编写派生数据集以Avro格式输出文件，而无需考虑模式。

## 模式的优点

正如上述所介绍的，Protocol Buffers、Thrift和Avro都使用了模式来描述二进制编码格式。它们的模式语言比XML模式或JSON模式简单得多，它支持更详细的验证规则（例如，“该字段的字符串值必须匹配此正则表达式”或“该字段的整数值必须介于0~100间”）。由于Protocol Buffers、Thrift和Avro的实现更简单，使用更简单，它们已经得到了非常广泛的编程语言支持。

不过，这些编码背后的思想绝不是什么新鲜事物。例如，它们与ASN.1有很多相似之处，ASN.1是在1984年首次被标准化的模式定义语言<sup>[27]</sup>。ASN.1被用来定义各种网络





协议，其二进制编码（DER）仍然被用于编码SSL证书（X.509）<sup>[28]</sup>。ASN.1支持使用标签号的模式演化，类似于Protocol Buffers和Thrift<sup>[29]</sup>。然而，它非常复杂，并且文档不尽如人意，所以ASN.1可能不是新型应用的最佳选择。

许多数据系统也实现了一些专有的二进制编码。例如，大多数关系数据库都有网络协议，可以通过该协议向数据库发送查询并获取响应。这些协议通常用于特定的数据库，并且数据库供应商提供驱动程序（例如，使用ODBC或JDBC API），把来自数据库的网络协议的响应解码为内存数据结构。

所以我们看到，尽管JSON、XML和CSV等文本数据格式非常普遍，但基于模式的二进制编码也是一个可行的选择。它们有许多不错的属性：

- 它们可以比各种“二进制JSON”变体更紧凑，可以省略编码数据中的字段名称。
- 模式是一种有价值的文档形式，因为模式是解码所必需的，所以可以确定它是最新的（而手动维护的文档可能很容易偏离现实）。
- 模式数据库允许在部署任何内容之前检查模式更改的向前和向后兼容性。
- 对于静态类型编程语言的用户来说，从模式生成代码的能力是有用的，它能够在编译时进行类型检查。

总之，通过演化支持与无模式/读时模式的JSON数据库相同的灵活性（参阅第2章“文档模型中的模式灵活性”），同时还提供了有关数据和工具方面更好的保障。

## 数据流模式

本章开始时，我们提到，每当将一些数据发送到非共享内存的另一个进程时，例如，当通过网络发送数据或者把它写入文件时，都需要将数据编码为字节序列。然后，讨论了用于执行此操作的各种不同编码技术。

而向前和向后的兼容性对于可演化性来说非常重要，通过允许独立升级系统的不同部分，而不必一次改变所有，使更改更为容易。兼容性是执行编码的一个进程和执行解码的另一个进程之间的关系。

这是一个相当抽象的想法，数据可以通过多种方式从一个进程流向另一个进程。谁编码数据？谁解码数据？在本章的其余部分，将探讨一些最常见的进程间数据流动的方式：



- 通过数据库（参阅本章后面的“基于数据库的数据流”部分）。
- 通过服务调用（参阅本章后面的“基于服务的数据流：REST和RPC”部分）。
- 通过异步消息传递（参阅本章后面的“基于消息传递的数据流”部分）。

## 基于数据库的数据流

在数据库中，写入数据库的进程对数据进行编码，而读取数据库的进程对数据进行解码。可能只有一个进程访问数据库，在这种情况下，reader只是同一进程的较新版本，此时，可以认为向数据库中存储内容，就是给未来的自己发送消息。

这种情况下，向后兼容性显然是必要的；否则未来的自己将无法解码以前写的东西。

一般而言，几个不同的进程同时访问数据库是很常见的。这些进程可能是几个不同的应用程序或服务，也可能只是同一服务的几个实例（为了可伸缩性或容错并行运行）。无论哪种情况，在应用程序正在改变的环境中，访问数据库的某些进程可能运行较新的代码，而某些进程可能运行较旧的代码。例如，因为当前正在滚动升级中部署新版本，所以某些实例已经更新，而其他实例尚未更新。

这意味着数据库中的值可以由较新版本的代码写入，然后由仍在运行的旧版本代码读取。因此，数据库通常也需要向前兼容。

然而，还有一个额外的障碍。假设在记录模式中添加了一个字段，并且较新的代码将该新字段的值写入数据库。随后，旧版本的代码（尚不知道该新字段）将读取、更新记录并将其写回。在这种情况下，理想的行为通常是旧代码保持新字段不变，即使它无法解释。

之前讨论的编码格式支持未知字段的保存，但是有时候还需要注意应用程序层面的影响，如图4-7所示。例如，如果将数据库值解码为应用程序中的模型对象，然后重新编码这些模型对象，则在该转换过程中可能会丢失未知字段。解决这个问题并不困难，重要的是首先要有这方面的意识。

### 不同的时间写入不同的值

数据库通常支持在任何时候更新任何值。这意味着在单个数据库中，可能有一些值是在5ms前写入的，而有些值是在5年前写入的。

部署新版本的应用程序（至少是服务器端应用程序）时，可能会在几分钟内用新版本完全替换旧版本。数据库内容的情况并不是这样：五年前的数据仍然采用原始编码，

除非已经明确地重写了它。这种现象有时被总结为数据比代码更长久。

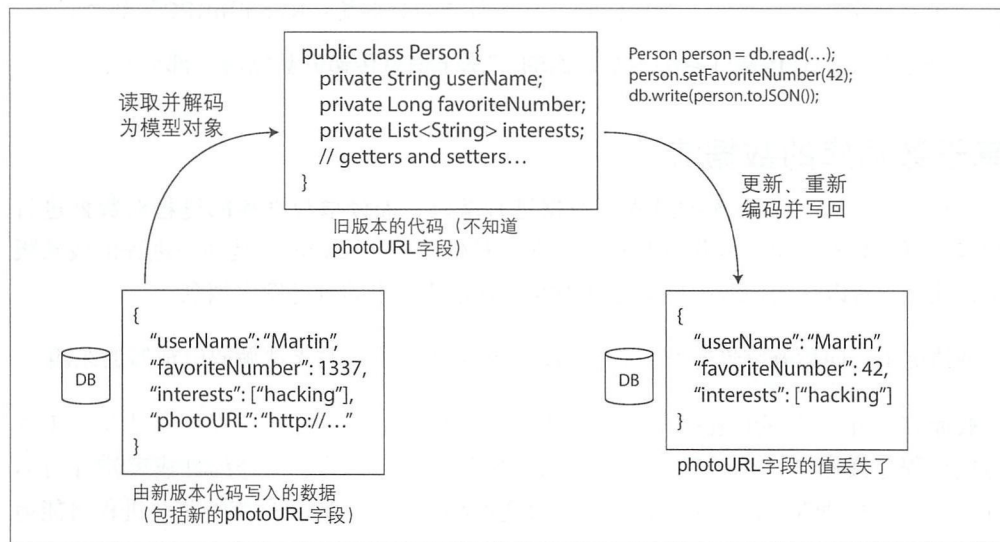


图4-7：当旧版本的应用程序更新新版本的程序所写入的数据时，需要小心，否则可能会丢失数据

将数据重写（或迁移）为新模式当然是可能的，但在大型数据集上执行此操作代价不菲，因此很多数据库都尽可能避免此操作。大多数关系数据库允许进行简单的模式更改，例如添加具有默认值为空的新列，而不重写现有数据<sup>注5</sup>。读取旧行时，数据库会为磁盘上编码数据缺失的所有列填充为空值。LinkedIn的文档数据库Espresso使用Avro进行存储，并支持Avro的模式演化规则<sup>[23]</sup>。

因此，模式演化支持整个数据库看起来像是采用单个模式编码，即使底层存储可能包含各个版本模式所编码的记录。

## 归档存储

或许你会不时地为数据库创建快照，例如用于备份或加载到数据仓库（参阅第3章的“数据仓库”部分）。在这种情况下，数据转储通常使用最新的模式进行编码，即使源数据库中的原始编码包含了不同时代的各种模式版本。由于无论如何都要复制数据，所以此时最好对数据副本进行统一的编码。

由于数据转储是一次写入的，而且以后不可改变，因此像Avro对象容器文件这样的格

注5：除了MySQL，MySQL经常会重写整个表，即使并非严格要求，如第2章“文档模型中的模式灵活性”所述。



式非常适合。这也是很好的机会，可以用分析友好的列存储（如Parquet）对数据进行编码（参阅第3章的“列压缩”部分）。

在第10章中，将详细讨论如何在归档存储中使用数据。

## 基于服务的数据流: REST和RPC

对于需要通过网络进行通信的进程，有多种不同的通信方式。最常见的是有两个角色：客户端和服务端。服务端通过网络公开API，客户端可以连接到服务端以向该API发出请求。服务端公开的API称为服务。

Web的工作方式是：客户端（Web浏览器）向Web服务端发出请求，发出GET请求来下载HTML、CSS、JavaScript、图像等，发出POST请求提交数据到服务端。API包含一组标准的协议和数据格式（HTTP、URL、SSL/TLS、HTML等）。因为Web浏览器、Web服务端和网站作者大多同意这些标准，所以可以使用任何浏览器访问任何网站（至少在理论上）。

Web浏览器不是唯一的客户端类型。例如，在移动设备或桌面计算机上运行的本地应用程序也可以向服务端发出网络请求，并且在Web浏览器内运行的客户端JavaScript应用程序可以使用XMLHttpRequest成为HTTP客户端（该技术被称为Ajax<sup>[30]</sup>）。在这种情况下，服务端的响应通常不是用于展示给用户的HTML，而是便于客户端应用程序代码进一步处理的编码数据（如JSON）。虽然HTTP可以用作传输协议，但是在顶层实现的API是特定于应用程序的，客户端和服务端需要就该API的细节达成一致。

此外，服务端本身可以是另一项服务的客户端（例如，典型的Web应用服务端作为数据库的客户端）。这种方法通常用于将大型应用程序按照功能区域分解为较小的服务，这样当一个服务需要另一个服务的某些功能或数据时，就会向另一个服务发出请求。这种构建应用程序的方式传统上被称为面向服务的体系结构（service-oriented architecture, SOA），最近则更名为微服务体系结构（microservices architecture）<sup>[31,32]</sup>。

在某些方面，服务类似于数据库：它们通常允许客户端提交和查询数据。然而，虽然数据库支持使用第2章中讨论的查询语言进行任意查询，但服务公开了特定于应用程序的API，它只允许由服务的业务逻辑（应用程序代码）预先确定的输入和输出<sup>[33]</sup>。此限制提供了一定程度的封装：服务可以对客户端可以做什么和不能做什么施加细粒度的限制。

面向服务/微服务体系结构的一个关键设计目标是，通过使服务可独立部署和演化，

让应用程序更易于更改和维护。例如，每个服务应该由一个团队拥有，该团队应该能够经常发布新版本的服务，而不必与其他团队协调。换句话说，应该期望新旧版本的服务器和客户端同时运行，因此服务器和客户端使用的数据编码必须在不同版本的服务API之间兼容，这正是在本章所讨论的内容。

## 网络服务

当HTTP被用作与服务通信的底层协议时，它被称为Web服务。这可能有点用词不当，因为Web服务不仅在Web上使用，而且在几个不同的上下文中使用。例如：

1. 运行在用户设备上的客户端应用程序（例如，移动设备上的本地应用程序，或使用Ajax的JavaScript Web应用程序），通过HTTP向服务发出请求。这些请求通常通过公共互联网进行。
2. 一种服务向同一组织拥有的另一项服务提出请求，这些服务通常位于同一数据中心内，作为面向服务/微型架构的一部分。支持这种用例的软件有时被称为中间件。
3. 一种服务向不同组织所拥有的服务提出请求，经常需通过互联网。这用于不同组织后端系统之间的数据交换。此类别包括由在线服务（如信用卡处理系统）提供的公共API，或用于共享访问用户数据的OAuth。

有两种流行的Web服务方法：REST和SOAP。它们在设计理念方面几乎是截然相反的，也往往是各自支持者之间激烈辩论的主题<sup>注6</sup>。

REST不是一种协议，而是一个基于HTTP原则的设计理念<sup>[34,35]</sup>。它强调简单的数据格式，使用URL来标识资源，并使用HTTP功能进行缓存控制、身份验证和内容类型协商。与SOAP相比，REST已经越来越受欢迎，至少在跨组织服务集成的背景下<sup>[36]</sup>，并经常与微服务相关联<sup>[31]</sup>。根据REST原则所设计的API称为RESTful。

相比之下，SOAP是一种基于XML的协议，用于发出网络API请求<sup>注7</sup>。虽然它最常用于HTTP，但其目的是独立于HTTP，并避免使用大多数HTTP功能。相反，它带有庞大而复杂的多种相关标准（Web服务框架，Web Service Framework，称为WS-\*）和新增的各种功能<sup>[37]</sup>。

---

注6： 即使在每个阵营内部也有很多争论。例如，HATEOAS (*hypermedia as the engine of application state*) 经常引发讨论<sup>[35]</sup>。

注7： 尽管首字母缩写词相似，但SOAP并不是SOA的必要条件。SOAP是一种特定的技术，而SOA是构建系统的一般方法。

SOAP Web服务的API使用被称为WSDL（Web Services Description Language，一种基于XML的语言）来描述。WSDL支持代码生成，客户端可以使用本地类和方法调用（编码为XML消息并由框架再次解码）来访问远程服务。这在静态类型编程语言中非常有用，但在动态类型编程语言中用处不大（参阅本章前面的“代码生成和动态类型化语言”部分）。

由于WSDL的设计目标不是人类可读的，而且SOAP消息通常过于复杂，无法手动构建，SOAP用户严重依赖工具支持、代码生成和IDE<sup>[38]</sup>。对于没有SOAP供应商支持的编程语言的用户来说，试图与SOAP服务集成非常困难。

尽管SOAP及其各种扩展表面上是标准化的，但是不同厂商的实现之间的互操作性往往存在一些问题<sup>[39]</sup>。由于所有这些原因，尽管它仍然被许多大型企业使用，但已经不再受到大多数小公司的青睐。

RESTful的API倾向于更简单的方法，通常涉及较少的代码生成和自动化工具。定义格式如OpenAPI，也称为Swagger<sup>[40]</sup>，可用于描述RESTful API并帮助生成文档。

## 远程过程调用（RPC）的问题

Web服务仅仅是通过网络发出API请求的一系列技术的最新体现，其中许多技术受到了大肆宣传，但也存在严重的问题。企业级JavaBeans（EJB）和Java的远程方法调用（RMI）仅限于Java。分布式组件对象模型（DCOM）仅限于Microsoft平台。公共对象请求代理体系结构（CORBA）过于复杂，不提供向前或向后兼容性<sup>[41]</sup>。

所有这些都是基于远程过程调用（Remote Procedure Call，RPC）的思想，该思想自20世纪70年代以来就一直存在<sup>[42]</sup>。RPC模型试图使向远程网络服务发出请求看起来与在同一进程中调用编程语言中的函数或方法相同（这种抽象称为位置透明）。虽然RPC起初看起来很方便，但这种方法在根本上是有缺陷的<sup>[43,44]</sup>。网络请求与本地函数调用非常不同：

- 本地函数调用是可预测的，并且成功或失败仅取决于控制的参数。网络请求是不可预测的：请求或响应可能由于网络问题而丢失，或者远程计算机可能速度慢或不可用，这些问题完全不在控制范围之内。网络问题很常见，因此必须有所准备，例如重试失败的请求。
- 本地函数调用要么返回一个结果，要么抛出一个异常，或者永远不会返回（因为进入无限循环或者进程崩溃）。网络请求有另一个可能的结果：由于超时，它返回时可能没有结果。在这种情况下，根本不知道发生了什么：如果没有收到来自远程服务的响应，无法知道请求是否成功。我们将在第8章更详细地讨论此类问题。



- 如果重试失败的网络请求，可能会发生请求实际上已经完成，只是响应丢失的情况。在这种情况下，重试将导致该操作被执行多次，除非在协议中建立重复数据消除（幂等性）机制。本地函数调用则没有这样问题。在第11章更详细地讨论幂等性。
- 每次调用本地函数时，通常需要大致相同的时间来执行。网络请求比函数调用要慢得多，而且其延迟也有很大的变化：情况好时，它可能会在不到1ms的时间内完成，但是当网络拥塞或者远程服务过载时，可能需要几秒钟的时间才能完成相同操作。
- 调用本地函数时，可以高效地将引用（指针）传递给本地内存中的对象。当发出网络请求时，所有这些参数都需要被编码成可以通过网络发送的字节序列。如果参数是像数字或字符串这样的基本类型，这没关系，但是对于较大的对象很快就会出现问題。
- 客户端和服务可以用不同的编程语言来实现，所以RPC框架必须将数据类型从一种语言转换成另一种语言。因为不是所有的语言都具有相同的类型，所以最终可能会丑陋。例如，回想一下JavaScript的数字大于 $2^{53}$ 的问题（参阅本章前面的“JSON、XML与二进制变体”部分）。用单一语言编写的单个进程中不存在此问题。

所有这些因素意味着，尝试使远程服务看起来像编程语言中的本地对象一样毫无意义，因为它们是根本不同的事情。REST的部分吸引力在于，它并不试图隐藏它是网络协议的事实（尽管这似乎并没有阻止人们在REST之上构建RPC库）。

## RPC的发展方向

虽然有这些问题，但是RPC并没有消失。在本章提到的所有编码的基础上构建了各种RPC框架：例如Thrift和Avro带有RPC支持，gRPC是使用Protocol Buffers的RPC实现，Finagle也使用Thrift，Rest.li使用HTTP上的JSON。

新一代的RPC框架更加明确了远程请求与本地函数调用不同的事实。例如，Finagle和Rest.li使用Futures（Promises）来封装可能失败的异步操作。Futures还简化了需要并行请求多项服务的情况，并将其结果合并<sup>[45]</sup>。gRPC支持流，其中调用不仅包括一个请求和一个响应，还包括一段时间内一系列的请求和响应<sup>[46]</sup>。

其中一些框架还提供了服务发现，即允许客户端查询在哪个IP地址和端口号上获得特定的服务。我们将在第6章“请求路由”回到该主题。

使用二进制编码格式的自定义RPC协议，可以实现比诸如REST上的JSON之类的通用

协议更好的性能。但是，RESTful API还有其他一些显著的优点：它有利于实验和调试（只需使用Web浏览器或命令行工具curl即可向它发出请求，而无需任何代码生成或软件安装），支持所有的主流编程语言和平台，并且有一个庞大的工具生态系统（服务器、缓存、负载均衡器、代理、防火墙、监控、调试工具、测试工具等）。

由于这些原因，REST似乎是公共API的主流风格。RPC框架主要侧重于同一组织内多项服务之间的请求，通常发生在同一数据中心内。

## RPC的数据编码和演化

对于演化性，重要的是可以独立地更改和部署RPC客户端和服务端。与基于数据库的数据流相比（如上一节所述），此处可以做一个简化的假设：假定所有的服务端都先被更新，其次是所有的客户端。因此，只需要在请求上具有向后兼容性，而在响应上具有向前兼容性。

RPC方案的向后和向前兼容性属性取决于它所使用的具体编码技术：

- Thrift、gRPC（Protocol Buffers）和Avro RPC可以根据各自编码格式的兼容性规则进行演化。
- 在SOAP中，请求和响应是用XML模式指定的。这些都是可以演化的，但有一些微妙的陷阱<sup>[47]</sup>。
- RESTful API通常使用JSON（没有正式指定的模式）用于响应，而请求则采用JSON或URI编码/表单编码的请求参数。为了保持兼容性，通常考虑的更改包括添加可选的请求参数和在响应中添加新的字段。

如果RPC经常用于跨组织边界的通信，则服务的兼容性会变得更加困难，服务的提供者经常无法控制其客户，也不能强制他们升级。因此，需要长期保持兼容性，也许是无限期的。如果不得不进行一些破坏兼容性的更改，则服务提供者往往会同时维护多个版本的服务API。

关于API版本管理应该如何工作（即客户端如何指示它想要使用哪个版本的API<sup>[48]</sup>）没有统一的方案。对于RESTful API，常用的方法是在URL或HTTP Accept头中使用版本号。对于使用API密钥来标识特定客户端的服务，另一种选择是将客户端请求的API版本存储在服务器上，并允许通过单独的管理接口更新该版本选项<sup>[49]</sup>。

## 基于消息传递的数据流

我们一直在研究从一个进程到另一个进程不同的数据流编码方式。到目前为止，已经

讨论了REST和RPC（其中一个进程通过网络向另一个进程发送请求，并期望尽快得到响应）以及数据库（其中一个进程写入编码数据，另一个进程在将来某个时刻再次读取该数据）。

在最后一节中，将简要介绍一下RPC和数据库之间的异步消息传递系统。它们与RPC的相似之处在于，客户端的请求（通常称为消息）以低延迟传递到另一个进程。它们与数据库的相似之处在于，不是通过直接的网络连接发送消息，而是通过称为消息代理（也称为消息队列，或面向消息的中间件）的中介发送的，该中介会暂存消息。

与直接RPC相比，使用消息代理有以下几个优点：

- 如果接收方不可用或过载，它可以充当缓冲区，从而提高系统的可靠性。
- 它可以自动将消息重新发送到崩溃的进程，从而防止消息丢失。
- 它避免了发送方需要知道接收方的IP地址和端口号（这在虚拟机经常容易起起停停的云部署中特别有用）。
- 它支持将一条消息发送给多个接收方。
- 它在逻辑上将发送方与接收方分离（发送方只是发布消息，并不关心谁使用它们）。

然而，与RPC的差异在于，消息传递通信通常是单向的：发送方通常不期望收到对其消息的回复。进程可能发送一个响应，但这通常是在一个独立的通道上完成的。这种通信模式是异步的：发送者不会等待消息被传递，而只是发送然后忘记它。

## 消息代理

过去，消息代理主要由TIBCO、IBM WebSphere和WebMethods等公司的商业软件所控制。最近，像RabbitMQ、ActiveMQ、HornetQ、NATS和Apache Kafka这样的开源实现已经流行起来。我们将在第11章对它们进行更详细的比较。

详细的传递语义因实现和配置而异，但通常情况下，消息代理的使用方式如下：一个进程向指定的队列或主题发送消息，并且代理确保消息被传递给队列或主题的一个或多个消费者或订阅者。在同一主题上可以有许多生产者与许多消费者。

主题只提供单向数据流。但是，消费者本身可能会将消息发布到另一个主题（因此可以将它们链接在一起，就像将在第11章中看到的那样），也可以发送到一个回复队列，该队列由原始消息发送者来消费（这样支持类似RPC的请求/响应数据流）。

消息代理通常不会强制任何特定的数据模型，消息只是包含一些元数据的字节序列，



因此可以使用任何编码格式。如果编码是向后和向前兼容的，则可以最大程度灵活地独立更改发布者和消费者，并以任意顺序部署他们。

如果消费者重新发布消息到另一个主题，则可能需要小心保留未知字段，以防止前面在数据库上下文中描述的问题（见图4-7）。

## 分布式Actor框架

Actor模型是用于单个进程中并发的编程模型。逻辑被封装在Actor中，而不是直接处理线程（以及竞争条件、锁定和死锁的相关问题）。每个Actor通常代表一个客户端或实体，它可能具有某些本地状态（不与其他任何Actor共享），并且它通过发送和接收异步消息与其他Actor通信。不保证消息传送：在某些错误情况下，消息将丢失。由于每个Actor一次只处理一条消息，因此不需要担心线程，每个Actor都可以由框架独立调度。

在分布式Actor框架中，这个编程模型被用来跨越多个节点来扩展应用程序。无论发送方和接收方是在同一个节点上还是在不同的节点上，都使用相同的消息传递机制。如果它们位于不同的节点上，则消息被透明地编码成字节序列，通过网络发送，并在另一端被解码。

相比RPC，位置透明性在Actor模型中更有效，因为Actor模型已经假定消息可能会丢失，即使在单个进程中也是如此。尽管网络上的延迟可能比同一个进程中的延迟更高，但是在使用Actor模型时，本地和远程通信之间根本上的不匹配所发生的概率更小。

分布式的Actor框架实质上是消息代理和Actor编程模型集成到单个框架中。但是，如果要对基于Actor的应用程序执行滚动升级，则仍需担心向前和向后兼容性问题，因为消息可能会从运行新版本的节点发送到运行旧版本的节点，反之亦然。

三种流行的分布式Actor框架处理消息编码的方式如下：

- 默认情况下，Akka使用Java的内置序列化，它不提供向前或向后兼容性。但是，可以用类似Protocol Buffers的东西替代它，从而获得滚动升级的能力<sup>[50]</sup>。
- 默认情况下，Orleans使用不支持滚动升级部署的自定义数据编码格式；要部署新版本的应用程序，需要建立一个新的集群，将流量从旧集群导入到新集群，然后关闭旧集群<sup>[51,52]</sup>。像Akka一样，也可以使用自定义序列化插件。
- 在Erlang OTP中，很难对记录模式进行更改（尽管系统具有许多为高可用性而设计的功能）。滚动升级在技术上是可能的，但要求仔细规划<sup>[53]</sup>。一些实验性的

- 新型映射数据类型（2014年在Erlang R17中引入的类似于JSON的结构）可能会使模式的更改在将来变得更容易<sup>[54]</sup>。

## 小结

本章，我们研究了将内存数据结构转换为网络或磁盘上字节流的多种方法。我们看到这些编码的细节不仅影响其效率，更重要的是还影响应用程序的体系结构和部署时的支持选项。

特别地，许多服务需要支持滚动升级，即每次将新版本的服务逐步部署到几个节点，而不是同时部署到所有节点。滚动升级允许在不停机的情况下发布新版本的服务（因此鼓励频繁地发布小版本而不是大版本），并降低部署风险（允许错误版本在影响大量用户之前检测并回滚）。这些特性非常有利于应用程序的演化和更改。

在滚动升级期间，或者由于各种其他原因，必须假设不同的节点正在运行应用代码的不同版本。因此，在系统内流动的所有数据都以提供向后兼容性（新代码可以读取旧数据）和向前兼容性（旧代码可以读取新数据）的方式进行编码显得非常重要。

本章还讨论了多种数据编码格式及其兼容性情况：

- 编程语言特定的编码仅限于某一种编程语言，往往无法提供向前和向后兼容性。
- JSON、XML和CSV等文本格式非常普遍，其兼容性取决于你如何使用他们。它们有可选的模式语言，这有时是有用的，有时却是一个障碍。这些格式对某些数据类型的支持有些模糊，必须小心处理数字和二进制字符串等问题。
- 像Thrift、Protocol Buffers和Avro这样的二进制的模式驱动格式，支持使用清晰定义的向前和向后兼容性语义进行紧凑、高效的编码。这些模式对于静态类型语言中的文档和代码生成非常有用。然而，它们有一个缺点，即只有在数据解码后才是人类可读的。

我们还讨论了数据流的几种模型，说明了数据编码在不同场景下非常重要：

- 数据库，其中写入数据库的进程对数据进行编码，而读取数据库的进程对数据进行解码。
- RPC和REST API，其中客户端对请求进行编码，服务器对请求进行解码并对响应进行编码，客户端最终对响应进行解码。

- 异步消息传递（使用消息代理或Actor），节点之间通过互相发送消息进行通信，消息由发送者编码并由接收者解码。

最后，我们可以得出这样的结论，只要稍加小心，向后/向前兼容性和滚动升级是完全可以实现的。预祝你的应用程序可以快速迭代，顺利部署。

## 参考文献

- [1] “Java Object Serialization Specification,” *docs.oracle.com*, 2010.
- [2] “Ruby 2.2.0 API Documentation,” *ruby-doc.org*, Dec 2014.
- [3] “The Python 3.4.3 Standard Library Reference Manual,” *docs.python.org*, February 2015.
- [4] “EsotericSoftware/kryo,” *github.com*, October 2014.
- [5] “CWE-502: Deserialization of Untrusted Data,” Common Weakness Enumeration, *cwe.mitre.org*, July 30, 2014.
- [6] Steve Breen: “What Do WebLogic, WebSphere, JBoss, Jenkins, OpenNMS, and Your Application Have in Common? This Vulnerability,” *foxglovesecurity.com*, November 6, 2015.
- [7] Patrick McKenzie: “What the Rails Security Issue Means for Your Startup,” *kalzumeus.com*, January 31, 2013.
- [8] Eishay Smith: “jvm-serializers wiki,” *github.com*, November 2014.
- [9] “XML Is a Poor Copy of S-Expressions,” *c2.com* wiki.
- [10] Matt Harris: “Snowflake: An Update and Some Very Important Information,” email to *Twitter Development Talk* mailing list, October 19, 2010.
- [11] Shudi (Sandy) Gao, C. M. Sperberg-McQueen, and Henry S. Thompson: “XML Schema 1.1,” W3C Recommendation, May 2001.
- [12] Francis Galiegue, Kris Zyp, and Gary Court: “JSON Schema,” IETF Internet-Draft, February 2013.
- [13] Yakov Shafranovich: “RFC 4180: Common Format and MIME Type for Comma-



Separated Values (CSV) Files,” October 2005.

[14] “MessagePack Specification,” *msgpack.org*.

[15] Mark Slee, Aditya Agarwal, and Marc Kwiatkowski: “Thrift: Scalable Cross-Language Services Implementation,” Facebook technical report, April 2007.

[16] “Protocol Buffers Developer Guide,” Google, Inc., *developers.google.com*.

[17] Igor Anishchenko: “Thrift vs Protocol Buffers vs Avro - Biased Comparison,” slideshare.net, September 17, 2012.

[18] “A Matrix of the Features Each Individual Language Library Supports,” *wiki.apache.org*.

[19] Martin Kleppmann: “Schema Evolution in Avro, Protocol Buffers and Thrift,” *martin.kleppmann.com*, December 5, 2012.

[20] “Apache Avro 1.7.7 Documentation,” *avro.apache.org*, July 2014.

[21] Doug Cutting, Chad Walters, Jim Kellerman, et al.: “[PROPOSAL] New Subproject: Avro,” email thread on *hadoop-general mailing list*, *mail-archives.apache.org*, April 2009.

[22] Tony Hoare: “Null References: The Billion Dollar Mistake,” at *QCon London*, March 2009.

[23] Aditya Auradkar and Tom Quiggle: “Introducing Espresso—LinkedIn’s Hot New Distributed Document Store,” *engineering.linkedin.com*, January 21, 2015.

[24] Jay Kreps: “Putting Apache Kafka to Use: A Practical Guide to Building a Stream Data Platform (Part 2),” *blog.confluent.io*, February 25, 2015.

[25] Gwen Shapira: “The Problem of Managing Schemas,” *radar.oreilly.com*, November 4, 2014.

[26] “Apache Pig 0.14.0 Documentation,” *pig.apache.org*, November 2014.

[27] John Larmouth: *ASN.1 Complete*. Morgan Kaufmann, 1999. ISBN: 978-0-122-33435-1.

[28] Russell Housley, Warwick Ford, Tim Polk, and David Solo: “RFC 2459: Internet

X.509 Public Key Infrastructure: Certificate and CRL Profile,” IETF Network Working Group, Standards Track, January 1999.

[29] Lev Walkin: “Question: Extensibility and Dropping Fields,” *lionet.info*, September 21, 2010.

[30] Jesse James Garrett: “Ajax: A New Approach to Web Applications,” *adaptivepath.com*, February 18, 2005.

[31] Sam Newman: *Building Microservices*. O’Reilly Media, 2015. ISBN: 978-1-491-95035-7.

[32] Chris Richardson: “Microservices: Decomposing Applications for Deployability and Scalability,” *infoq.com*, May 25, 2014.

[33] Pat Helland: “Data on the Outside Versus Data on the Inside,” at *2nd Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2005.

[34] Roy Thomas Fielding: “Architectural Styles and the Design of Network-Based Software Architectures,” PhD Thesis, University of California, Irvine, 2000.

[35] Roy Thomas Fielding: “REST APIs Must Be Hypertext-Driven,” *roy.gbiv.com*, October 20 2008.

[36] “REST in Peace, SOAP,” *royal.pingdom.com*, October 15, 2010.

[37] “Web Services Standards as of Q1 2007,” *innoc.com*, February 2007.

[38] Pete Lacey: “The S Stands for Simple,” *harmful.cat-v.org*, November 15, 2006.

[39] Stefan Tilkov: “Interview: Pete Lacey Criticizes Web Services,” *infoq.com*, December 12, 2006.

[40] “OpenAPI Specification (fka Swagger RESTful API Documentation Specification) Version 2.0,” *swagger.io*, September 8, 2014.

[41] Michi Henning: “The Rise and Fall of CORBA,” *ACM Queue*, volume 4, number 5, pages 28–34, June 2006. doi:10.1145/1142031.1142044.

[42] Andrew D. Birrell and Bruce Jay Nelson: “Implementing Remote Procedure Calls,” *ACM Transactions on Computer Systems (TOCS)*, volume 2, number 1, pages

39-59, February 1984. doi:10.1145/2080.357392.

[43] Jim Waldo, Geoff Wyant, Ann Wollrath, and Sam Kendall: “A Note on Distributed Computing,” Sun Microsystems Laboratories, Inc., Technical Report TR-94-29, November 1994.

[44] Steve Vinoski: “Convenience over Correctness,” *IEEE Internet Computing*, volume 12, number 4, pages 89–92, July 2008. doi:10.1109/MIC.2008.75.

[45] Marius Eriksen: “Your Server as a Function,” at *7th Workshop on Programming Languages and Operating Systems (PLOS)*, November 2013. doi:10.1145/2525528.2525538.

[46] “grpc-common Documentation,” Google, Inc., *github.com*, February 2015.

[47] Aditya Narayan and Irina Singh: “Designing and Versioning Compatible Web Services,” *ibm.com*, March 28, 2007.

[48] Troy Hunt: “Your API Versioning Is Wrong, Which Is Why I Decided to Do It 3 Different Wrong Ways,” *troyhunt.com*, February 10, 2014.

[49] “API Upgrades,” Stripe, Inc., April 2015.

[50] Jonas Bonér: “Upgrade in an Akka Cluster,” email to *akka-user mailing list*, *grokbase.com*, August 28, 2013.

[51] Philip A. Bernstein, Sergey Bykov, Alan Geller, et al.: “Orleans: Distributed Virtual Actors for Programmability and Scalability,” Microsoft Research Technical Report MSR-TR-2014-41, March 2014.

[52] “Microsoft Project Orleans Documentation,” Microsoft Research, *dotnet.github.io*, 2015.

[53] David Mercer, Sean Hinde, Yinso Chen, and Richard A O’Keefe: “beginner: Updating Data Structures,” email thread on *erlang-questions* mailing list, *erlang.com*, October 29, 2007.

[54] Fred Hebert: “Postscript: Maps,” *learnyousomeerlang.com*, April 9, 2014.



## 第二部分

# 分布式数据系统

成功的技术应首先处理好现实问题，因为现实无法被愚弄。

——Richard Feynman，罗杰斯委员会报告（1986）

本书第一部分主要讨论了单台机器存储系统设计的主要技术。在第二部分，我们将继续向前迈进，当需要多台机器提供数据存储和检索服务时，又会有哪些挑战和方案呢？

主要出于以下目的，我们需要在多台机器上分布数据：

### 扩展性

当数据量或者读写负载巨大，严重超出了单台机器的处理上限，需要将负载分散到多台机器上。

### 容错与高可用性

当单台机器（或者多台，以及网络甚至整个数据中心）出现故障，还希望应用系统可以继续工作，这时需要采用多台机器提供冗余。这样某些组件失效之后，冗余组件可以迅速接管。

### 延迟考虑

- 如果客户遍布世界各地，通常需要考虑在全球范围内部署服务，以方便用户就近访问最近数据中心所提供的服务，从而避免数据请求跨越了半个地球才能到达目标。

## 系统扩展能力

当负载增加需要更强的处理能力时，最简单的办法就是购买更强大的机器（有时称为垂直扩展）。由一个操作系统管理更多的CPU，内存和磁盘，通过高速内部总线使每个CPU都可以访问所有的存储器或磁盘。在这样一个共享内存架构中，所有这些组件的集合可看作一台大机器<sup>[1]注1</sup>。

共享内存架构的问题在于，成本增长过快甚至超过了线性：即如果把一台机器内的CPU数量增加一倍，内存扩容一倍，磁盘容量加大一倍，则最终总成本增加不止一倍。并且由于性能瓶颈因素，这样一台机器尽管拥有了两倍的硬件指标但却不一定能处理两倍的负载。

共享内存架构能够提供有限的容错能力，例如高端的服务器可以热插拔很多组件（在不关闭机器的情况下更换磁盘，内存模块，甚至是CPU）。但很显然，它仍局限于某个特定的地理位置，无法提供异地容错能力。

另一种方法是共享磁盘架构，它拥有多台服务器，每个服务器各自拥有独立的CPU和内存，然后将数据存储在可共享访问的磁盘阵列上，服务器与磁盘阵列之间往往通过高速网络连接<sup>注2</sup>。这种架构多适用于数据仓库等负载，然而通常由于资源竞争以及锁的开销等限制了其进一步的扩展能力<sup>[2]</sup>。

## 无共享结构

相比之下，无共享架构（也称为水平扩展）<sup>[3]</sup>则获得了很大的关注度。当采用这种架构时，运行数据库软件的机器或者虚拟机称为节点。每个节点独立使用本地的CPU，内存和磁盘。节点之间的所有协调通信等任务全部运行在传统网络（以太网）之上且核心逻辑主要依靠软件来实现。

无共享系统不需要专门的硬件，具有较高的性价比。它可以跨多个地理区域分发数据，从而减少用户的访问延迟，甚至当整个数据中心发生灾难时仍能继续工作。通过

---

注1：常见的高端服务器中，尽管每个CPU可以访问所有内存，但由于物理总线连接原因，CPU与某些内存更为靠近（即非对称内存访问，或称NUMA<sup>[11]</sup>）。为了更加高效地利用这种NUMA架构，应用处理最好根据NUMA拓扑结构进行划分，使CPU主要访问本地内存。这意味着即使在一台机器内，仍然需要某种分区。

注2：主要指网络附加存储（Network Attached Storage, NAS）或者存储局域网（Storage Area Network, SAN）。

云计算虚拟机的部署方式，即便是没有Google级别规模的小公司，也可以轻松拥有跨区域的分布式架构和服务能力。

本部分内容将重点放在无共享体系架构上，并不是因为它一定是所有应用的最佳选择，而是因为它需要应用开发者更多的关注和深入理解。例如把数据分布在多节点上，就需要了解在这样一个分布式系统下，背后的权衡设计和隐含限制，数据库并不能魔法般地把所有复杂性都屏蔽起来。

虽然分布式无共享体系架构具有很多优点，但也会给应用程序带来更多的复杂性，有时甚至会限制实际可用的数据模型。例如在某些极端情况下，一个简单的单线程程序可能比一个拥有100多个CPU核的集群性能更好<sup>[4]</sup>。而另一方面，无共享系统也可以做到性能非常强大。接下来的几章里我们将详细讨论数据分布时所面临的主要问题。

## 复制与分区

将数据分布在多节点时有两种常见的方式：

### 复制

在多个节点上保存相同数据的副本，每个副本具体的存储位置可能不尽相同。复制方法可以提供冗余：如果某些节点发生不可用，则可以通过其他节点继续提供数据访问服务。复制也可以帮助提高系统性能。我们在第5章将主要讨论复制技术。

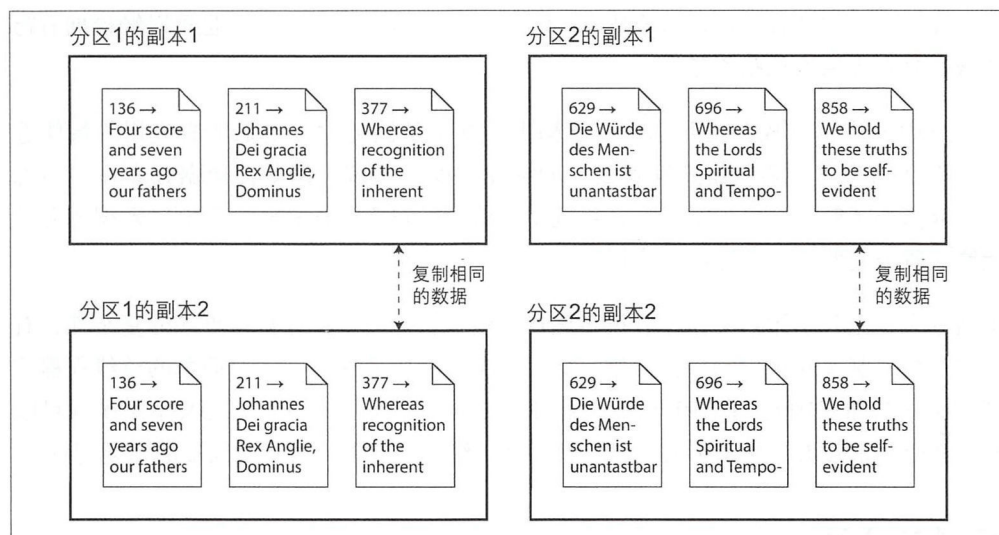
### 分区

将一个大块头的数据库拆分成多个较小的子集即分区，不同的分区分配给不同的节点（也称为分片）。我们在第6章主要介绍分区技术。

这些是不同的数据分布机制，然而它们经常被放在一起组合使用。参见图II-1的示例。

在了解以上概念之后，我们会讨论分布式环境中错综复杂的权衡之道，很可能你会在实际设计系统时不得不面对这些艰难选择。我们将在第7章介绍事务，帮助理解数据系统中各种可能出错的情况以及处理方法，而第8章和第9章将深入分析分布式系统内在的局限性，之后结束本部分。





图II-1：数据库有两个分区，每个分区各有两份副本

在之后本书的第三部分，我们将讨论如何把多个（可能每一个都是分布式）数据存储组件集成到一个更大的系统中，以满足更复杂的应用需求。但在那之前，我们首先来谈谈分布式数据系统。

## 参考文献

- [1] Ulrich Drepper: “What Every Programmer Should Know About Memory,” *akkadia.org*, November 21, 2007.
- [2] Ben Stopford: “Shared Nothing vs. Shared Disk Architectures: An Independent View,” *benstopford.com*, November 24, 2009.
- [3] Michael Stonebraker: “The Case for Shared Nothing,” *IEEE Database Engineering Bulletin*, volume 9, number 1, pages 4–9, March 1986.
- [4] Frank McSherry, Michael Isard, and Derek G. Murray: “Scalability! But at What COST?,” at *15th USENIX Workshop on Hot Topics in Operating Systems (HotOS)*, May 2015.



# 数据复制

一个可能出错的事物与一个不可能出错的事物之间的主要区别是，当一个不可能出错的事物出错了，通常也就意味着不可修复。

——Douglas Adams, 《基本无害》（1992）

复制主要指通过互联网络在多台机器上保存相同数据的副本。正如第二部分开头所介绍的，通过数据复制方案，人们通常希望达到以下目的：

- 使数据在地理位置上更接近用户，从而降低访问延迟。
- 当部分组件出现故障，系统依然可以继续工作，从而提高可用性。
- 扩展至多台机器以同时提供数据访问服务，从而提高读吞吐量。

本章我们将假设数据规模比较小，集群的每一台机器都可以保存数据集的完整副本。在接下来的第6章中，我们放宽这一假设，讨论单台机器无法容纳整个数据集的情况（即必须分区）。在后面的章节中，我们还将讨论复制过程中可能出现的各种故障，以及该如何处理这些故障。

如果复制的数据一成不变，那么复制就非常容易：只需将数据复制到每个节点，一次即可搞定。然而所有的技术挑战都在于处理那些持续更改的数据，而这正是本章讨论的核心。我们将讨论三种流行的复制数据变化的方法：主从复制、多主节点复制和无主节点复制。几乎所有的分布式数据库都使用上述方法中的某一种，而三种方法各有优缺点，我们稍后会详细解读。

复制技术存在许多需要折中考虑的地方，例如采用同步复制还是异步复制，以及如何



处理失败的副本等。数据库通常采用可配置选项来调整这些处理策略，虽然在处理细节方面因数据库实现而异，但存在一些通用的一般性原则。本章我们还将讨论不同选项可能出现的后果。

数据库复制其实是个很古老的话题。因为网络的基本约束条件没有发生本质的改变，可以说自1970年所研究的基本复制原则，时至今日也没有发生太大的变化<sup>[1]</sup>。然而，除了学术研究，实践中很多开发人员仍然假定数据库只运行在单节点上，分布式数据库成为主流也只是最近发生的事情。由于许多应用开发人员在这方面经验还略显不足，对诸如“最终一致性”等问题存在一些误解。因此，在“复制滞后问题”中，我们会详细讨论最终一致性，包括读自己的写和单调读等。

## 主节点与从节点

每个保存数据库完整数据集的节点称之为副本。当有了多副本，不可避免地会引入一个问题：如何确保所有副本之间的数据是一致的？

对于每一笔数据写入，所有副本都需要随之更新；否则，某些副本将出现不一致。最常见的解决方案是基于主节点的复制（也称为主动/被动，或主从复制），如图5-1所示。主从复制的工作原理如下：

1. 指定某一个副本为主副本（或称为主节点）。当客户写数据库时，必须将写请求首先发送给主副本，主副本首先将新数据写入本地存储。
2. 其他副本则全部称为从副本（或称为从节点）<sup>注1</sup>。主副本把新数据写入本地存储后，然后将数据更改作为复制的日志或更改流发送给所有从副本。每个从副本获得更改日志之后将其应用到本地，且严格保持与主副本相同的写入顺序。
3. 客户端从数据库中读数据时，可以在主副本或者从副本上执行查询。再次强调，只有主副本才可以接受写请求；从客户端的角度来看，从副本都是只读的。

许多关系型数据库都内置支持主从复制，例如PostgreSQL（9.0版本以后）、MySQL、Oracle Data Guard<sup>[2]</sup>和SQL Server的AlwaysOn Availability Groups<sup>[3]</sup>。而一些非关系数据库如MongoDB、RethinkDB和Espresso<sup>[4]</sup>也支持主从复制。另外，主从复制技术也不仅限于数据库，还广泛用于分布式消息队列如Kafka<sup>[5]</sup>和RabbitMQ<sup>[6]</sup>，以及一些网络文件系统和复制块设备（如DRBD）。

---

注1： 存在不同的冷、热、冷备服务器的定义。例如PostgreSQL中，热备是指负责接受客户端读的副本，而warm standby则负责处理主副本角色的变化事件，但不负责客户端的任何查询。由于并不妨碍总体理解，本书忽略这些差异。

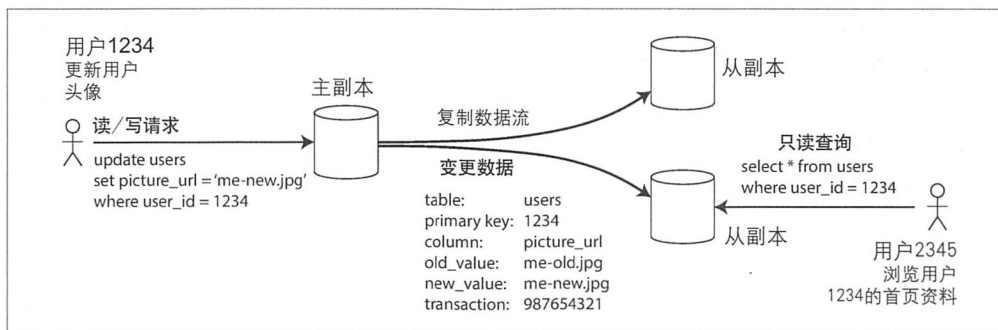


图5-1：主从复制系统

## 同步复制与异步复制

复制非常重要的一个设计选项是同步复制还是异步复制。对于关系数据库系统，同步或异步通常是一个可配置的选项；而其他系统则可能是硬性指定或者只能二选一。

结合图5-1的例子，网站用户需要更新首页的头像图片。其基本流程是，客户将更新请求发送给主节点，主节点接收到请求，接下来将数据更新转发给从节点。最后，由主节点来通知客户更新完成。

图5-2则进一步描述了系统各个模块间的通信情况，包括客户端，主节点和两个从节点。时间从左到右。请求或响应标记为粗箭头。

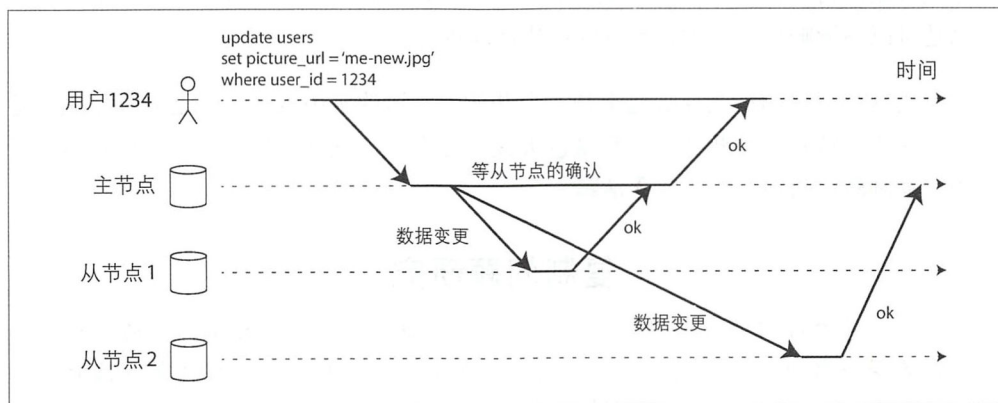


图5-2：主从复制，包括一个同步的从节点和一个异步的从节点

图5-2中，从节点1的复制是同步的，即主节点需等待直到从节点1确认完成了写入，然后才会向用户报告完成，并且将最新的写入对其他客户端可见。而从节点2的复制是异步的：主节点发送完消息之后立即返回，不用等待从节点2的完成确认。

图5-2中，从节点2在接收复制日志之前有一段很长的延迟。通常情况下，复制速度会非常快，例如多数数据库系统可以在一秒之内完成所有从节点的更新。但是，系统其实并没有保证一定会在多长时间内完成复制。有些情况下，从节点可能落后主节点几分钟甚至更长时间，例如，由于从节点刚从故障中恢复，或者系统已经接近最大设计上限，或者节点之间的网络出现问题。

同步复制的优点是，一旦向用户确认，从节点可以明确保证完成了与主节点的更新同步，数据已经处于最新版本。万一主节点发生故障，总是可以在从节点继续访问最新数据。缺点则是，如果同步的从节点无法完成确认（例如由于从节点发生崩溃，或者网络故障，或任何其他原因），写入就不能视为成功。主节点会阻塞其后所有的写操作，直到同步副本确认完成。

因此，把所有从节点都配置为同步复制有些不切实际。因为这样的话，任何一个同步节点的中断都会导致整个系统更新停滞不前。实践中，如果数据库启用了同步复制，通常意味着其中某一个从节点是同步的，而其他节点则是异步模式。万一同步的从节点变得不可用或性能下降，则将另一个异步的从节点提升为同步模式。这样可以保证至少有两个节点（即主节点和一个同步从节点）拥有最新的数据副本。这种配置有时也称为半同步<sup>[7]</sup>。

主从复制还经常会被配置为全异步模式。此时如果主节点发生失败且不可恢复，则所有尚未复制到从节点的写请求都会丢失。这意味着即使向客户端确认了写操作，却无法保证数据的持久化。但全异步配置的优点则是，不管从节点上数据多么滞后，主节点总是可以继续响应写请求，系统的吞吐性能更好。

异步模式这种弱化的持久性听起来是一个非常不靠谱的折中设计，但是异步复制还是被广泛使用，特别是那些从节点数量巨大或者分布于广域地理环境。我们将在本章后面的“复制滞后问题”继续这个话题。

## 复制问题研究

主节点发生故障时异步复制系统可能会丢失数据，这是一个非常严重的问题，因此在保证数据不丢失的前提下，研究人员尝试了各种办法来提高复制性能与系统可用性。例如，链式复制<sup>[8,9]</sup>是同步复制的一种变体，已经在一些系统（如Microsoft Azure存储<sup>[10,11]</sup>）中得以实现。

多副本一致性与共识之间有着密切的联系（即让多个节点对数据状态达成一致），我们将在第9章详细探讨这一点。本章主要集中于数据库实践中常用的、相对简单的复制技术。



## 配置新的从节点

当如果出现以下情况时，如需要增加副本数以提高容错能力，或者替换失败的副本，就需要考虑增加新的从节点。但如何确保新的从节点和主节点保持数据一致呢？

简单地将数据文件从一个节点复制到另一个节点通常是不够的。主要是因为客户端仍在不断向数据库写入新数据，数据始终处于不断变化之中，因此常规的文件拷贝方式将会导致不同节点上呈现出不同时间点的数据，这不是我们所期待的。

或许应该考虑锁定数据库（使其不可写）来使磁盘上的文件保持一致，但这会违反高可用的设计目标。好在我们可以做到在不停机、数据服务不中断的前提下完成从节点的设置。逻辑上的主要操作步骤如下：

1. 在某个时间点对主节点的数据副本产生一个一致性快照，这样避免长时间锁定整个数据库。目前大多数数据库都支持此功能，快照也是系统备份所必需的。而在某些情况下，可能需要第三方工具，如MySQL的innobackupex<sup>[12]</sup>。
2. 将此快照拷贝到新的从节点。
3. 从节点连接到主节点并请求快照点之后所发生的数据更改日志。因为在第一步创建快照时，快照与系统复制日志的某个确定位置相关联，这个位置信息在不同的系统有不同的称呼，如PostgreSQL将其称为“log sequence number”（日志序列号），而MySQL将其称为“binlog coordinates”。
4. 获得日志之后，从节点来应用这些快照点之后所有数据变更，这个过程称之为追赶。接下来，它可以继续处理主节点上新的数据变化。并重复步骤1~步骤4。

建立新的从副本具体操作步骤可能因数据库系统而异。某些系统中，这个过程是全自动化的，而其他系统中所涉及的步骤、流程可能会比较复杂，甚至需要管理员手动介入。

## 处理节点失效

系统中的任何节点都可能因故障或者计划内的维护（例如重启节点以安装内核安全补丁）而导致中断甚至停机。如果能够在不停机的情况下重启某个节点，这会对运维带来巨大的便利。我们的目标是，尽管个别节点会出现中断，但要保持系统总体的持续运行，并尽可能减小节点中断带来的影响。

那么如何通过主从复制技术来实现系统高可用呢？

## 从节点失效：追赶式恢复

从节点的本地磁盘上都保存了副本收到的数据变更日志。如果从节点发生崩溃，然后顺利重启，或者主从节点之间的网络发生暂时中断（闪断），则恢复比较容易，根据副本的复制日志，从节点可以知道在发生故障之前所处理的最后一笔事务，然后连接到主节点，并请求自那笔事务之后中断期间内所有的数据变更。在收到这些数据变更日志之后，将其应用到本地来追赶主节点。之后就和正常情况一样持续接收来自主节点数据流的变化。

## 主节点失效：节点切换

处理主节点故障的情况则比较棘手：选择某个从节点将其提升为主节点；客户端也需要更新，这样之后的写请求会发送给新的主节点，然后其他从节点要接受来自新的主节点上的变更数据，这一过程称之为切换。

故障切换可以手动进行，例如通知管理员主节点发生失效，采取必要的步骤来创建新的主节点；或者以自动方式进行。自动切换的步骤通常如下：

1. **确认主节点失效。**有很多种出错可能性，例如由于系统崩溃，停电，网络问题等。没有万无一失的方法能够确切地检测到究竟问题出在哪里，所以大多数系统都采用了基于超时的机制：节点间频繁地互相发生发送心跳存活消息，如果发现某一个节点在一段比较长时间（例如30s）没有响应，即认为该节点发生失效（如果主节点在计划内出于维护目的而故意下线，则不在此讨论范围）。
2. **选举新的主节点。**可以通过选举的方式（超过多数的节点达成共识）来选举新的主节点，或者由之前选定的某控制节点来指定新的主节点。候选节点最好与原主节点的数据差异最小，这样可以最小化数据丢失的风险。让所有节点同意新的主节点是个典型的共识问题，会在第9章详细讨论。
3. **重新配置系统使新主节点生效。**客户端现在需要将写请求发送给新的主节点（细节将在第6章的“请求路由”中讨论）。如果原主节点之后重新上线，可能仍然自认为是主节点，而没有意识到其他节点已经达成共识迫使其下台。这时系统要确保原主节点降级为从节点，并认可新的主节点。

然而，上述切换过程依然充满了很多变数：

- 如果使用了异步复制，且失效之前，新的主节点并未收到原主节点的所有数据；在选举之后，原主节点很快又重新上线并加入到集群，接下来的写操作会发生什么？新的主节点很可能会收到冲突的写请求，这是因为原主节点未意识的角色变化，还会尝试同步其他从节点，但其中的一个现在已经接管成为现任主节点。

常见的解决方案是，原主节点上未完成复制的写请求就此丢弃，但这可能会违背数据更新持久化的承诺。

- 如果在数据库之外有其他系统依赖于数据库的内容并在一起协同使用，丢弃数据的方案就特别危险。例如，在GitHub的一个事故中<sup>[13]</sup>，某个数据并非完全同步的MySQL从节点被提升为主副本，数据库使用了自增计数器将主键分配给新建的行，但是因为新的主节点计数器落后于原主节点（即二者并非完全同步），它重新使用了已被原主节点分配出去的某些主键，而恰好这些主键已被外部Redis所引用，结果出现MySQL和Redis之间的不一致，最后导致了某些私有数据被错误地泄露给了其他用户。
- 在某些故障情况下（参见第8章），可能会发生两个节点同时都自认为是主节点。这种情况被称为脑裂，它非常危险：两个主节点都可能接受写请求，并且没有很好解决冲突的办法（参阅本章后面的“多主节点复制技术”），最后数据可能会丢失或者破坏。作为一种安全应急方案，有些系统会采取措施来强制关闭其中一个节点<sup>注2</sup>。然而，如果设计或者实现考虑不周，可能会出现两个节点都被关闭的情况<sup>[14]</sup>。
- 如何设置合适的超时来检测主节点失效呢？主节点失效后，超时时间设置得越长也意味着总体恢复时间就越长。但如果超时设置太短，可能会导致很多不必要的切换。例如，突发的负载峰值会导致节点的响应时间变长甚至超时，或者由于网络故障导致延迟增加。如果系统此时已经处于高负载压力或网络已经出现严重拥塞，不必要的切换操作只会使总体情况变得更糟。

坦白讲，对于这些问题没有简单的解决方案。因此，即使系统可能支持自动故障切换，有些运维团队仍然更愿意以手动方式来控制整个切换过程。

上述这些问题，包括节点失效、网络不可靠、副本一致性、持久性、可用性与延迟之间各种细微的权衡，实际上正是分布式系统核心的基本问题。在第8章和第9章中，我们还会进一步讨论。

## 复制日志的实现

主从复制技术到底是如何工作呢？实践中有多种不同的实现方法，此处我们逐一做些介绍。

---

注2：该方法也被称为fencing（隔离）。在第8章中的“主节点与锁”会有更多细节讨论。



## 基于语句的复制

最简单的情况，主节点记录所执行的每个写请求（操作语句）并将该操作语句作为日志发送给从节点。对于关系数据库，这意味着每个INSERT、UPDATE或DELETE语句都会转发给从节点，并且每个从节点都会分析并执行这些SQL语句，如同它们是来自客户端那样。

听起来很合理也不复杂，但这种复制方式有一些不适用的场景：

- 任何调用非确定性函数的语句，如NOW()获取当前时间，或RAND()获取一个随机数等，可能会在不同的副本上产生不同的值。
- 如果语句中使用了自增列，或者依赖于数据库的现有数据（例如，UPDATE ... WHERE <某些条件>），则所有副本必须按照完全相同的顺序执行，否则可能会带来不同的结果。进而，如果有多个同时并发执行的事务时，会有很大的限制。
- 有副作用的语句（例如，触发器、存储过程、用户定义的函数等），可能会在每个副本上产生不同的副作用。

有可能采取一些特殊措施来解决这些问题，例如，主节点可以在记录操作语句时将非确定性函数替换为执行之后的确定的结果，这样所有节点直接使用相同的结果值。但是，这里面存在太多边界条件需要考虑，因此目前通常首选的是其他复制实现方案。

MySQL 5.1版本之前采用基于操作语句的复制。现在由于逻辑紧凑，依然在用，但是默认情况下，如果语句中存在一些不确定性操作，则MySQL会切换到基于行的复制（稍后讨论）。VoltDB使用基于语句的复制，它通过事务级别的确定性来保证复制的安全<sup>[15]</sup>。

## 基于预写日志（WAL）传输

在第3章中，我们讨论了存储引擎的磁盘数据结构，通常每个写操作都是以追加写的方式写入到日志中：

- 对于日志结构存储引擎（参阅第3章的“SSTables和LSM-trees”），日志是主要的存储方式。日志段在后台压缩并支持垃圾回收。
- 对于采用覆盖写磁盘的Btree（参阅第3章的“B-tree”）结构，每次修改会预先写入日志，如系统发生崩溃，通过索引更新的方式迅速恢复到此前一致状态。

不管哪种情况，所有对数据库写入的字节序列都被记入日志。因此可以使用完全相同

的日志在另一个节点上构建副本：除了将日志写入磁盘之外，主节点还可以通过网络将其发送给从节点。

从节点收到日志进行处理，建立和主节点内容完全相同的数据副本。

PostgreSQL、Oracle以及其他系统<sup>[16]</sup>等支持这种复制方式。其主要缺点是日志描述的数据结果非常底层：一个WAL包含了哪些磁盘块的哪些字节发生改变，诸如此类的细节。这使得复制方案和存储引擎紧密耦合。如果数据库的存储格式从一个版本改为另一个版本，那么系统通常无法支持主从节点上运行不同版本的软件。

看起来这似乎只是个有关实现方面的小细节，但可能对运营产生巨大的影响。如果复制协议允许从节点的软件版本比主节点更新，则可以实现数据库软件的不停机升级：首先升级从节点，然后执行主节点切换，使升级后的从节点成为新的主节点。相反，复制协议如果要求版本必须严格一致（例如WAL传输），那么就势必以停机为代价。

## 基于行的逻辑日志复制

另一种方法是复制和存储引擎采用不同的日志格式，这样复制与存储逻辑剥离。这种复制日志称为逻辑日志，以区分物理存储引擎的数据表示。

关系数据库的逻辑日志通常是指一系列记录来描述数据表行级别的写请求：

- 对于行插入，日志包含所有相关列的新值。
- 对于行删除，日志里有足够的信息来唯一标识已删除的行，通常是靠主键，但如果表上没有定义主键，就需要记录所有列的旧值。
- 对于行更新，日志包含足够的信息来唯一标识更新的行，以及所有列的新值（或至少包含所有已更新列的新值）。

如果一条事务涉及多行的修改，则会产生多个这样的日志记录，并在后面跟着一条记录，指出该事务已经提交。MySQL的二进制日志binlog（当配置为基于行的复制时）使用该方式<sup>[17]</sup>。

由于逻辑日志与存储引擎逻辑解耦，因此可以更容易地保持向后兼容，从而使主从节点能够运行不同版本的软件甚至是不同的存储引擎。

对于外部应用程序来说，逻辑日志格式也更容易解析。如果要将数据库的内容发送到外部系统（如用于离线分析的数据仓库），或构建自定义索引和缓存<sup>[18]</sup>等，基于逻辑日志的复制更有优势。该技术也被称为变更数据捕获，我们将在第11章中继续讨论。

## 基于触发器的复制

到目前为止所描述的复制方法都是由数据库系统来实现的，不涉及任何应用程序代码。通常这是大家所渴望的，不过，在某些情况下，我们可能需要更高的灵活性。例如，只想复制数据的一部分，或者想从一种数据库复制到另一种数据库，或者需要订制、管理冲突解决逻辑（参阅本章后面的“处理写冲突”），则需要将复制控制交给应用程序层。

有一些工具，例如Oracle GoldenGate<sup>[19]</sup>，可以通过读取数据库日志让应用程序获取数据变更。另一种方法则是借助许多关系数据库都支持的功能：触发器和存储过程。

触发器支持注册自己的应用层代码，使得当数据库系统发生数据更改（写事务）时自动执行上述自定义代码。通过触发器技术，可以将数据更改记录到一个单独的表中，然后外部处理逻辑访问该表，实施必要的自定义应用层逻辑，例如将数据更改复制到另一个系统。Oracle的Databus<sup>[20]</sup>和Postgres的Bucardo 就是这种技术的典型代表。

基于触发器的复制通常比其他复制方式开销更高，也比数据库内置复制更容易出错，或者暴露一些限制。然而，其高度灵活性仍有用武之地。

## 复制滞后问题

容忍节点故障只是使用复制其中的一个原因。正如第二部分开头所介绍的，其他原因包括可扩展性（采用多节点来处理更多的请求）和低延迟（将副本部署在地理上距离用户更近的地方）。

主从复制要求所有写请求都经由主节点，而任何副本只能接受只读查询。对于读操作密集的负载（如Web），这是一个不错的选择：创建多个从副本，将读请求分发给这些从副本，从而减轻主节点负载并允许读请求就近满足。

在这种扩展体系下，只需添加更多的从副本，就可以提高读请求的服务吞吐量。但是，这种方法实际上只能用于异步复制，如果试图同步复制所有的从副本，则单个节点故障或网络中断将使整个系统无法写入。而且节点越多，发生故障的概率越高，所以完全同步的配置现实中反而非常不可靠。

不幸的是，如果一个应用正好从一个异步的从节点读取数据，而该副本落后于主节点，则应用可能会读到过期的信息。这会导致数据库中出现明显的不一致：由于并非所有的写入都反映在从副本上，如果同时对主节点和从节点发起相同的查询，可能会得到不同的结果。这种不一致只是一个暂时的状态，如果停止写数据库，经过一



段时间之后，从节点最终会赶上并与主节点保持一致。这种效应也被称为最终一致性<sup>[22,23]注3</sup>。

“最终”一词有些含糊不清，总的来说，副本落后的程度理论上并没有上限。正常情况下，主节点和从节点上完成写操作之间的时间延迟（复制滞后）可能不足1秒，这样的滞后，在实践中通常不会导致太大影响。但是，如果系统已接近设计上限，或者网络存在问题，则滞后可能轻松增加到几秒甚至几分钟不等。

当滞后时间太长时，导致的不一致性不仅仅是一个理论存在的问题，而是个实实在在的现实问题。在本节中，我们将重点介绍三个复制滞后可能出现的问题，并给出相应的解决思路。

### 读自己的写

许多应用让用户提交一些数据，接下来查看他们自己所提交的内容。例如客户数据库中的记录，亦或者是讨论主题的评论等。提交新数据须发送到主节点，但是当用户读取数据时，数据可能来自从节点。这对于读密集和偶尔写入的负载是个非常合适的方案。

然而对于异步复制存在这样一个问题，如图5-3所示，用户在写入不久即查看数据，则新数据可能尚未到达从节点。对用户来讲，看起来似乎是刚刚提交的数据丢失了，显然用户不会高兴。

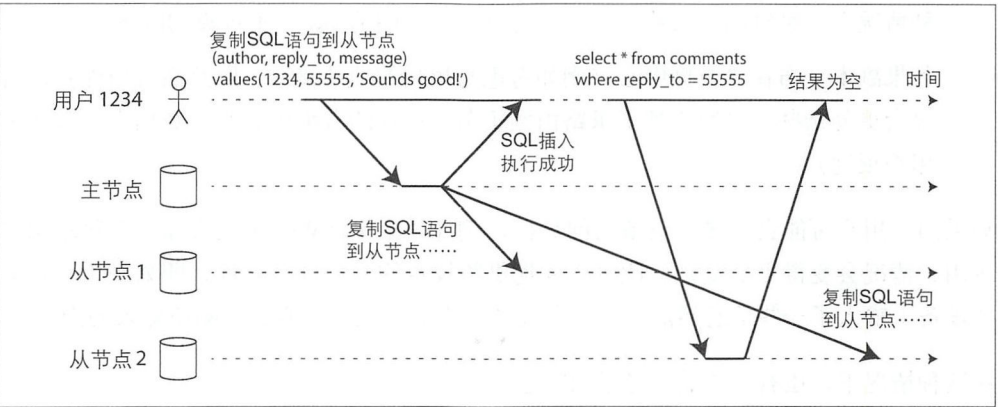


图5-3：用户发起写请求，然后在滞后的副本上读数据。此时需要read-after-write 一致性

注3： 最终一致性一词最早由Douglas Terry<sup>[24]</sup>发明，后来经由Werner Volgels<sup>[22]</sup>而普及，并成为很多NoSQL项目的标志性口号。但事实上，采用异步复制的关系数据库同样具备该特性。

对于这种情况，我们需要“写后读一致性”，也称为读写一致性<sup>[24]</sup>。该机制保证如果用户重新加载页面，他们总能看到自己最近提交的更新。但对其他用户则没有任何保证，这些用户的更新可能会在稍后才能刷新看到。

基于主从复制的系统该如何实现写后读一致性呢？有多种可行的方案，以下列举一二：

- 如果用户访问可能会被修改的内容，从主节点读取；否则，在从节点读取。这背后就要求有一些方法在实际执行查询之前，就已经知道内容是否可能会被修改。例如，社交网络上的用户首页信息通常只能由所有者编辑，而其他人无法编辑。因此，这就形成一个简单的规则：总是从主节点读取用户自己的首页配置文件，而在从节点读取其他用户的配置文件。
- 如果应用的大部分内容都可能被所有用户修改，那么上述方法将不太有效，它会导致大部分内容都必须经由主节点，这就丧失了读操作的扩展性。此时需要其他方案来判断是否从主节点读取。例如，跟踪最近更新的时间，如果更新后一分钟之内，则总是在主节点读取；并监控从节点的复制滞后程度，避免从那些滞后时间超过一分钟的从节点读取。
- 客户端还可以记住最近更新时的时间戳，并附带在读请求中，据此信息，系统可以确保对该用户提供读服务时都应该至少包含了该时间戳的更新。如果不够新，要么交由另一个副本来处理，要么等待直到副本接收到了最近的更新。时间戳可以是逻辑时间戳（例如用来指示写入顺序的日志序列号）或实际系统时钟（在这种情况下，时钟同步又称为一个关键点，请参阅第8章“不可靠的时钟”）。
- 如果副本分布在多数据中心（例如考虑与用户的地理接近，以及高可用性），情况会更复杂些。必须先把请求路由到主节点所在的数据中心（该数据中心可能离用户很远）。

如果同一用户可能会从多个设备访问数据，例如一个桌面Web浏览器和一个移动端的应用，情况会变得更加复杂。此时，要提供跨设备的写后读一致性，即如果用户在某个设备上输入了一些信息然后在另一台设备上查看，也应该看到刚刚所输入的内容。

在这种情况下，还有一些需要考虑的问题：

- 记住用户上次更新时间戳的方法实现起来会比较困难，因为在一台设备上运行的代码完全无法知道在其他设备上发生了什么。此时，元数据必须做到全局共享。
- 如果副本分布在多数据中心，无法保证来自不同设备的连接经过路由之后都到达

同一个数据中心。例如，用户的台式计算机使用了家庭宽带连接，而移动设备则使用蜂窝数据网络，不同设备的网络连接线路可能完全不同。如果方案要求必须从主节点读取，则首先需要想办法确保将来自不同设备的请求路由到同一个数据中心。

## 单调读

在前面异步复制读异常的第二个例子里，出现了用户数据向后回滚的奇怪情况。

假定用户从不同副本进行了多次读取，如图5-4所示，用户刷新一个网页，读请求可能被随机路由到某个从节点。用户2345先后在两个从节点上执行了两次完全相同的查询（先是少量滞后的节点，然后是滞后很大的从节点），则很有可能出现以下情况。第一个查询返回了最近用户1234所添加的评论，但第二个查询因为滞后的原因，还没有收到更新因而返回结果是空。实际上，第二个查询结果代表了更早时间点的状态。如果第一个查询没有返回任何内容，用户2345并不知道用户1234最近的评论，情况还不算太糟糕，但当用户2345看到了用户1234的评论之后，紧接着评论又消失了，他就会感觉很困惑。

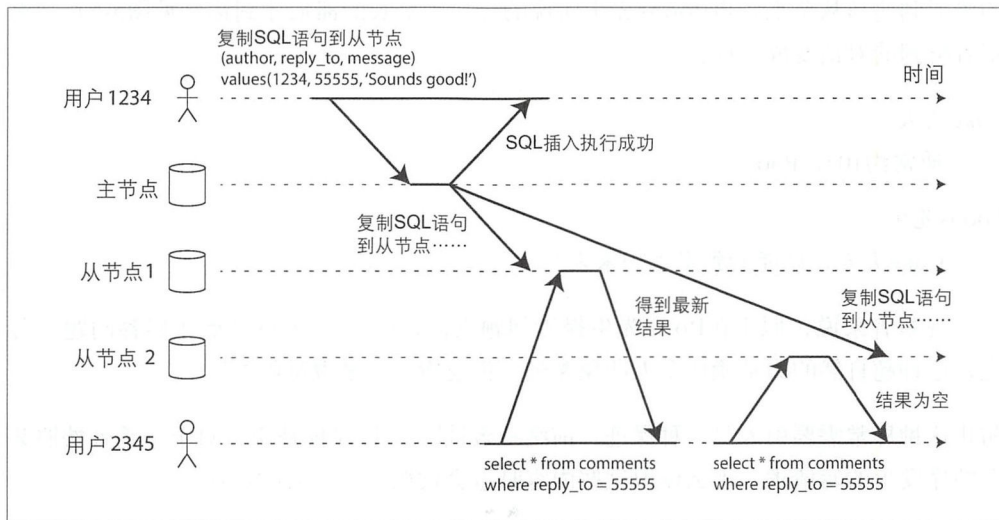


图5-4：用户看到了最新内容之后又读到了过期的内容，好像时间被回拨，此时需要单调读一致性

单调读<sup>[23]</sup>一致性可以确保不会发生这种异常。这是一个比强一致性弱，但比最终一致性强的保证。当读取数据时，单调读保证，如果某个用户依次进行多次读取，则他绝不会看到回滚现象，即在读取较新值之后又发生读旧值的情况。



实现单调读的一种方式，确保每个用户总是从固定的同一副本执行读取（而不同的用户可以从不同的副本读取）。例如，基于用户ID的哈希的方法而不是随机选择副本。但如果该副本发生失效，则用户的查询必须重新路由到另一个副本。

## 前缀一致读

第三个由于复制滞后导致因果反常的例子。例如Poons先生与Cake夫人之间的以下对话：

Poons先生

Cake夫人，您能看到多远的未来？

Cake夫人

通常约10s，Poons先生。

这两句话之间存在因果关系：Cake夫人首先是听到了Poons先生的问题，然后再去回答该问题。

现在，想象第三个人正在通过从节点收听上述对话。Cake夫人所说的话经历了短暂的滞后到达该节点，但Poons先生所说的经历了更长的滞后才到达（见图5-5）。观察者听到的对话变成这样：

Cake夫人

通常约10s，Poons先生。

Poons先生

Cake夫人，您能看到多远的未来？

对于观察者来说，似乎在Poon先生提出问题之前，Cake夫人就开始了回答问题。首先，这种超自然的力量确认令人印象深刻，但逻辑上却是混乱的<sup>[25]</sup>。

防止这种异常需要引入另一种保证：前缀一致读<sup>[23]</sup>。该保证是说，对于一系列按照某个顺序发生的写请求，那么读取这些内容时也会按照当时写入的顺序。

这是分区（分片）数据库中出现的一个特殊问题，细节将在第6章中讨论。如果数据库总是以相同的顺序写入，则读取总是看到一致的序列，不会发生这种反常。然而，在许多分布式数据库中，不同的分区独立运行，因此不存在全局写入顺序。这就导致当用户从数据库中读数据时，可能会看到数据库的某部分旧值和另一部分新值。

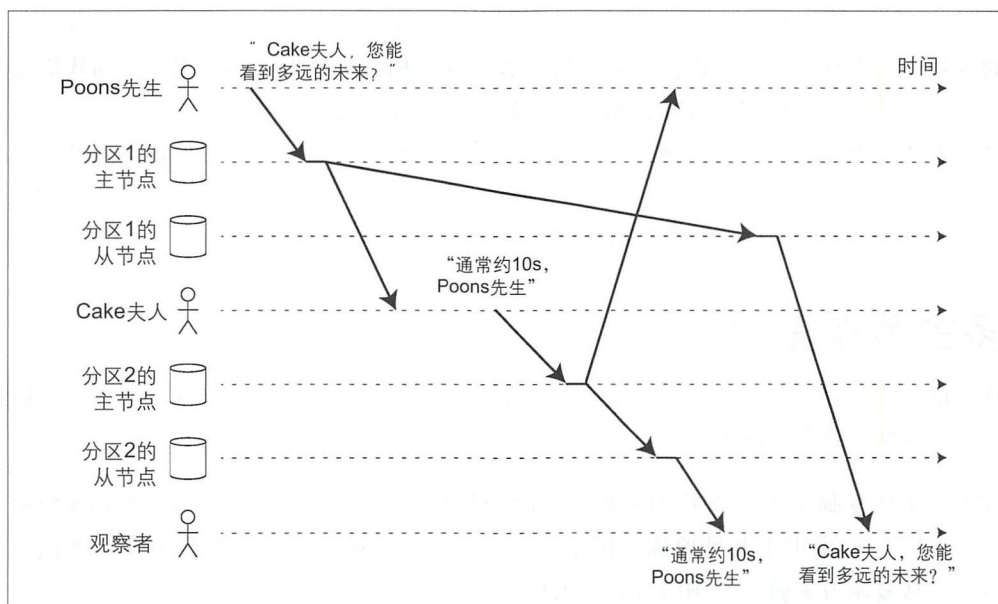


图5-5：分区数据经多副本复制后出现了不同程度的滞后，导致观察者先看到果后看到因

一个解决方案是确保任何具有因果顺序关系的写入都交给一个分区来完成，但该方案真实实现效率会大打折扣。现在有一些新的算法来显式地追踪事件因果关系，在本章稍后的“Happened-before关系与并发”会继续该问题的探讨。

## 复制滞后的解决方案

使用最终一致性系统时，最好事先就思考这样的问题：如果复制延迟增加到几分钟甚至几小时，那么应用层的行为会是什么样子？如果答案是“没问题”，那没得说。但是，如果带来糟糕的用户体验，那么在设计系统时，就要考虑提供一个更强的一致性保证，比如写后读；如果系统设计时假定是同步复制，但最终它事实上成为了异步复制，就可能会导致灾难性后果。

正如前面所讨论的，在应用层可以提供比底层数据库更强有力的保证。例如只在主节点上进行特定类型的读取，而代价则是，应用层代码中处理这些问题通常会非常复杂，且容易出错。

如果应用程序开发人员不必担心这么多底层的复制问题，而是假定数据库在“做正确的事情”，情况就变得很简单。而这也是事务存在的原因，事务是数据库提供更强保证的一种方式。

单节点上支持事务已经非常成熟。然而，在转向分布式数据库（即支持复制和分区）的过程中，有许多系统却选择放弃支持事务，并声称事务在性能与可用性方面代价过高，然后断言在可扩展的分布式系统中最终的一致性是无法避免的终极选择。关于这样的表述，首先它有一定道理，但情况远不是所说的那么简单，我们将在本书其余部分展开讨论，尝试形成一个更为深入的观点。例如在第7章和第9章将理解事务，然后在第三部分再介绍其他一些替代机制。

## 多主节点复制

到目前为止，我们只考虑了单个主节点的主从复制架构。主从复制方法较为常见，但也存在其他一些有趣的方案。

首先，主从复制存在一个明显的缺点：系统只有一个主节点，而所有写入都必须经由主节点<sup>注4</sup>。如果由于某种原因，例如与主节点之间的网络中断而导致主节点无法连接，主从复制方案就会影响所有的写入操作。

对主从复制模型进行自然的扩展，则可以配置多个主节点，每个主节点都可以接受写操作，后面复制的流程类似：处理写的每个主节点都必须将该数据更改转发到所有其他节点。这就是多主节点（也称为主-主，或主动/主动）复制。此时，每个主节点还同时扮演其他主节点的从节点。

## 适用场景

在一个数据中心内部使用多主节点基本没有太大意义，其复杂性已经超过所能带来的好处。但是，在以下场景这种配置则是合理的。

### 多数据中心

为了容忍整个数据中心级别故障或者更接近用户，可以把数据库的副本横跨多个数据中心。而如果使用常规的基于主从的复制模型，主节点势必只能放在其中的某一个数据中心，而所有写请求都必须经过该数据中心。

有了多主节点复制模型，则可以在每个数据中心都配置主节点，如图5-6所示的基本架构。在每个数据中心内，采用常规的主从复制方案；而在数据中心之间，由各个数据中心的主节点来负责同其他数据中心的主节点进行数据的交换、更新。

---

注4：如果数据库还采用了分区，不同的分区可能将主副本放在不同的节点上，但对于给定的某分区，则只有一个主节点。



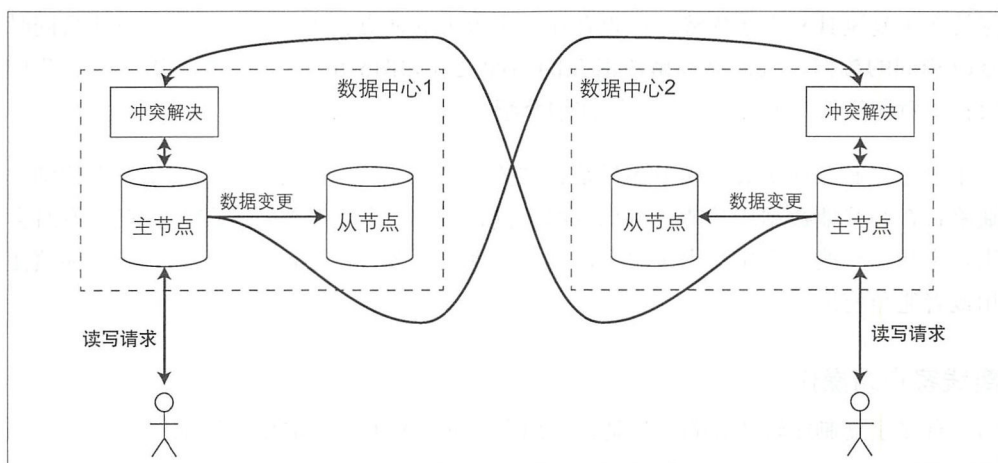


图5-6：横跨多个数据中心的 多主节点复制模型

可以对比一下在多个数据中心环境下，部署单主节点的主从复制方案与多主复制方案之间的差异：

### 性能

对于主从复制，每个写请求都必须经由广域网传送至主节点所在的数据中心。这会大大增加写入延迟，并基本偏离了采用多个数据中心的初衷（即就近访问）。而在多主节点模型中，每个写操作都可以在本地数据中心快速响应，然后采用异步复制方式将变化同步到其他数据中心。因此，对上层应用有效屏蔽了数据中心之间的网络延迟，使得终端用户所体验到的性能更好。

### 容忍数据中心失效

对于主从复制，如果主节点所在的数据中心发生故障，必须切换至另一个数据中心，将其中的一个从节点被提升为主节点。在多主节点模型中，每个数据中心则可以独立于其他数据中心继续运行，发生故障的数据中心在恢复之后更新到最新状态。

### 容忍网络问题

数据中心之间的通信通常经由广域网，它往往不如数据中心内的本地网络可靠。对于主从复制模型，由于写请求是同步操作，对数据中心之间的网络性能和稳定性等更加依赖。多主节点模型则通常采用异步复制，可以更好地容忍此类问题，例如临时网络闪断不会妨碍写请求最终成功。

有些数据库已内嵌支持了多主复制，但有些则借助外部工具来实现，例如MySQL的Tungsten Replicator<sup>[26]</sup>，PostgreSQL的BDR<sup>[27]</sup>以及Oracle的GoldenGate<sup>[19]</sup>。

尽管多主复制具有上述优势，但也存在一个很大的缺点：不同的数据中心可能会同时修改相同的数据，因而必须解决潜在的写冲突（如图5-6中的“冲突解决”）。我们会在本章稍后的“处理写入冲突”详细介绍。

由于多主复制在许多数据库中还只是新增的高级功能，所以可能存在配置方面的细小缺陷；在与其他数据库功能（例如自增主键，触发器和完整性约束等）交互时有时会出现意想不到的副作用。出于这个原因，一些人认为多主复制比较危险，应该谨慎使用或者避免使用<sup>[28]</sup>。

## 离线客户端操作

另一种多主复制比较适合的场景是，应用在与网络断开后还需要继续工作。

比如手机，笔记本电脑和其他设备上的日历应用程序。无论设备当前是否联网，都需要能够随时查看当前的会议安排（对应于读请求）或者添加新的会议（对应于写请求）。在离线状态下进行的任何更改，会在下次设备上线时，与服务器以及其他设备同步。

这种情况下，每个设备都有一个充当主节点的本地数据库（用来接受写请求），然后在所有设备之间采用异步方式同步这些多主节点上的副本，同步滞后可能是几小时或者数天，具体时间取决于设备何时可以再次联网。

从架构层面来看，上述设置基本上等同于数据中心之间的多主复制，只不过是极端情况，即一个设备就是数据中心，而且它们之间的网络连接非常不可靠。多个设备同步日历的例子表明，多主节点可以得到想要的结果，但中间过程依然有很多的未知数。

有一些工具可以使多主配置更为容易，如CouchDB就是为这种操作模式而设计的<sup>[29]</sup>。

## 协作编辑

实时协作编辑应用程序允许多个用户同时编辑文档。例如，Etherpad<sup>[30]</sup>和Google Docs<sup>[31]</sup>允许多人同时编辑文本文档或电子表格（算法简要会在本章后面的“自动冲突解决”中讨论）。

我们通常不会将协作编辑完全等价于数据库复制问题，但二者确实有很多相似之处。当一个用户编辑文档时，所做的更改会立即应用到本地副本（Web浏览器或客户端应用程序），然后异步复制到服务器以及编辑同一文档的其他用户。

如果要确保不会发生编辑冲突，则应用程序必须先将文档锁定，然后才能对其进行编

辑。如果另一个用户想要编辑同一个文档，首先必须等到第一个用户提交修改并释放锁。这种协作模式相当于主从复制模型下在主节点上执行事务操作。

为了加快协作编辑的效率，可编辑的粒度需要非常小。例如，单个按键甚至是全程无锁。然而另一方面，也会面临所有多主复制都存在的挑战，即如何解决冲突<sup>[32]</sup>。

## 处理写冲突

多主复制的最大问题是可能发生写冲突，这意味着必须有方案来解决冲突。

例如，两个用户同时编辑Wiki页面，如图5-7所示。用户1将页面的标题从A更改为B，与此同时用户2却将标题从A改为C。每个用户的更改都顺利地提交到本地主节点。但是，当更改被异步复制到对方时，却发现存在冲突<sup>[33]</sup>。注意，正常情况下的主从复制则不会出现这种情况。

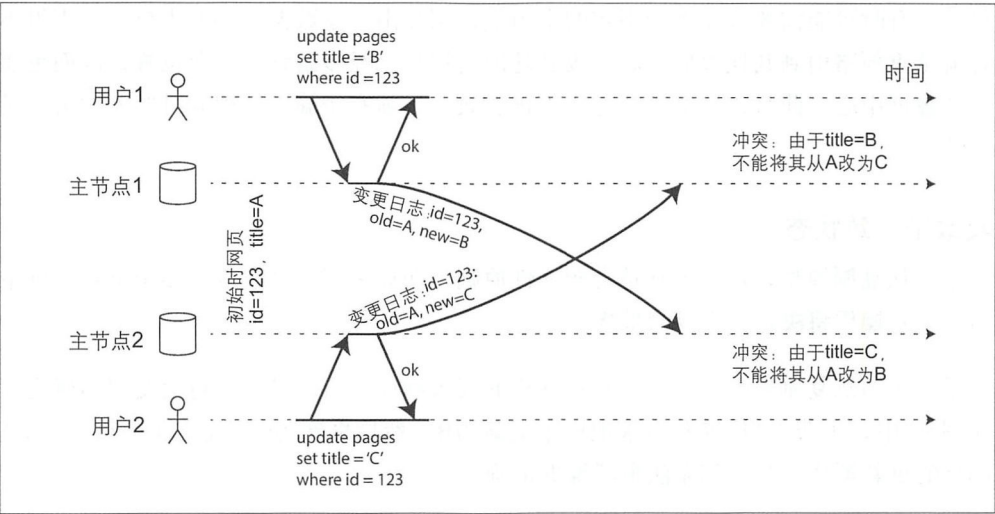


图5-7：两个主节点同时更新同一条记录而导致写冲突

## 同步与异步冲突检测

如果是主从复制数据库，第二个写请求要么会被阻塞直到第一个写完成，要么被中止（用户必须重试）。然而在多主节点的复制模型下，这两个写请求都是成功的，并且只能在稍后的时间点上才能异步检测到冲突，那时再要求用户层来解决冲突为时已晚。

理论上，也可以做到同步冲突检测，即等待写请求完成对所有副本的同步，然后再通知用户写入成功。但是，这样做将会失去多主节点的主要优势：允许每个主节点独立



接受写请求。如果确实想要同步方式冲突检测，或许应该考虑采用单主节点的主从复制模型。

## 避免冲突

处理冲突最理想的策略是避免发生冲突，即如果应用层可以保证对特定记录的写请求总是通过同一个主节点，这样就不会发生写冲突。现实中，由于不少多主节点复制模型所实现的冲突解决方案存在瑕疵，因此，避免冲突反而成为大家普遍推荐的首选方案<sup>[34]</sup>。

例如，一个应用系统中，用户需要更新自己的数据，那么我们确保特定用户的更新请求总是路由到特定的数据中心，并在该数据中心的主节点上进行读/写。不同的用户则可能对应不同的主数据中心（例如根据用户的地理位置来选择）。从用户的角度来看，这基本等价于主从复制模型。

但是，有时可能需要改变事先指定的主节点，例如由于该数据中心发生故障，不得不将流量重新路由到其他数据中心，或者是因为用户已经漫游到另一个位置，因而更靠近新数据中心。此时，冲突避免方式不再有效，必须有措施来处理同时写入冲突的可能性。

## 收敛于一致状态

对于主从复制模型，数据更新符合顺序性原则，即如果同一个字段有多个更新，则最后一个写操作将决定该字段的最终值。

对于多主节点复制模型，由于不存在这样的写入顺序，所以最终值也会变得不确定。在图5-7中，主节点1接受到请求把标题更新为B，然后更新为C；而在主节点2，则是相反的更新顺序。两者都无法辩驳谁更正确。

如果每个副本都只是按照它所看到写入的顺序执行，那么数据库最终将处于不一致状态。例如主节点1看到最终值C，而主节点2看到的是B，这绝对是不可接受的，所有的复制模型至少应该确保数据在所有副本中最终状态一定是一致的。因此，数据库必须以一种收敛趋同的方式来解决冲突，这也意味着当所有更改最终被复制、同步之后，所有副本的最终值是相同的。

实现收敛的冲突解决有以下可能的方式：

- 给每个写入分配唯一的ID，例如，一个时间戳，一个足够长的随机数，一个UUID或者一个基于键-值的哈希，挑选最高ID的写入作为胜利者，并将其他写入

丢弃。如果基于时间戳，这种技术被称为最后写入者获胜。虽然这种方法很流行，但是很容易造成数据丢失<sup>[35]</sup>。我们将在本章最后部分来详细解释。

- 为每个副本分配一个唯一的ID，并制定规则，例如序号高的副本写入始终优先于序号低的副本。这种方法也可能会导致数据丢失。
- 以某种方式将这些值合并在一起。例如，按字母顺序排序，然后拼接在一起（图5-7中，合并的标题可能类似于“B / C”）。
- 利用预定义好的格式来记录和保留冲突相关的所有信息，然后依靠应用层的逻辑，事后解决冲突（可能会提示用户）。

## 自定义冲突解决逻辑

解决冲突最合适的方式可能还是依靠应用层，所以大多数多主节点复制模型都有工具来让用户编写应用代码来解决冲突。可以在写入时或在读取时执行这些代码逻辑：

### 在写入时执行

只要数据库系统在复制变更日志时检测到冲突，就会调用应用层的冲突处理程序。例如，Bucardo支持编写一段Perl代码。这个处理程序通常不能在线提示用户，而只能在后台运行，这样速度更快。

### 在读取时执行

当检测到冲突时，所有冲突写入值都会暂时保存下来。下一次读取数据时，会将数据的多个版本读返回给应用层。应用层可能会提示用户或自动解决冲突，并将最后的结果返回到数据库。CouchDB采用了这样的处理方式。

注意，冲突解决通常用于单个行或文档，而不是整个事务<sup>[36]</sup>。因此，如果有一个原子事务包含多个不同写请求（如第7章），每个写请求仍然是分开考虑来解决冲突。

## 什么是冲突？

有些冲突是显而易见的。在图5-7的例子中，两个写操作同时修改同一个记录中的同一个字段，并将其设置为不同的值。毫无疑问，这就是一个冲突。

而其他类型的冲突可能会非常微妙，更难以发现。例如一个会议室预订系统，它主要记录哪个房间由哪个人在哪个时间段所预订。这个应用程序需要确保每个房间只能有一组人同时预定（即不得有相同房间的重复预订）。如果为同一个房间创建两个不同的预订，可能会发生冲突。尽管应用在预订时会检查房间是否可用，但如果两个预订是在两个不同的主节点上进行，则还是存在冲突的可能。

## 自动冲突解决

冲突解决的规则可能会变得越来越复杂，且自定义代码很容易出错。亚马逊是一个经常被引用的反面例子：有一段时间，购物的冲突解决逻辑依靠用户的购物车页面，后者保存了所有的物品，但顾客有时候会发现之前已经被拿掉的商品，再次出现在他们的购物车中<sup>[37]</sup>。

有一些有意思的研究尝试自动解决并发修改所引起的冲突。下面这些方法值得一看：

1. 无冲突的复制数据类型（Conflict-free Replicated Datatypes, CRDT）<sup>[32,38]</sup>。CRDT是可以由多个用户同时编辑的数据结构，包括map、ordered list、计数器等，并且以内置的合理方式自动地解决冲突。一些CRDT已经在Riak 2.0中得以具体实现<sup>[39,40]</sup>。
2. 可合并的持久数据结构（Mergeable persistent data）<sup>[41]</sup>。它跟踪变更历史，类似于Git版本控制系统，并提出三向合并功能（three-way merge function, CRDT采用双向合并）。
3. 操作转换（Operational transformation<sup>[42]</sup>）。它是Etherpad<sup>[30]</sup>和Google Docs<sup>[31]</sup>等协作编辑应用背后的冲突解决算法。专为可同时编辑的有序列表而设计，如文本文档的字符列表。

这些算法总体来讲还处于早期阶段，但将来它们可能会被整合到更多的数据系统中。这些自动冲突解决方案可以使主复制模型更简单、更容易被应用程序来集成。

很遗憾，这里没有现成的答案。通过接下来的几章，我们将对这个问题进行深入的剖析和讲解。我们将在第7章看到更多的冲突示例，在第12章中讨论检测 and 解决冲突的可扩展方法。

## 拓扑结构

复制的拓扑结构描述了写请求从一个节点的传播到其他节点的通信路径。如果有两个主节点，如图5-7所示，则只存在一个合理的拓扑结构：主节点1必须把所有的写同步到主节点2，反之亦然。但如果存在两个以上的主节点，则会有多个可能的同步拓扑结构，如图5-8所示。



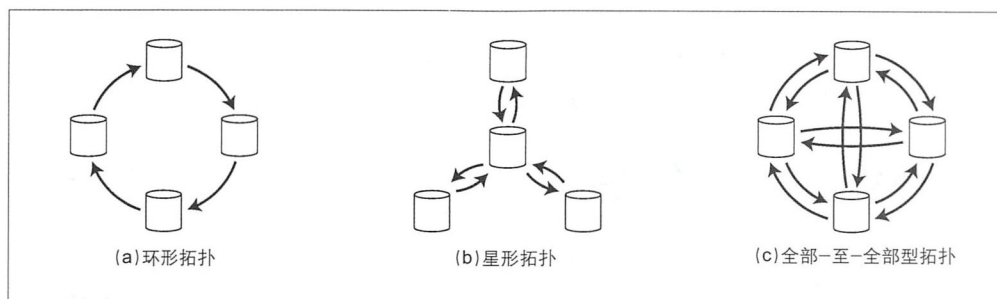


图5-8：多主节点模型的三种拓扑结构示例

最常见的拓扑结构是全部-至-全部，见图5-8 (c)，每个主节点将其写入同步到其他所有主节点。而其他一些拓扑结构也有普遍使用，例如，默认情况下MySQL只支持环形拓扑结构<sup>[34]</sup>，其中的每个节点接收来自前序节点的写入，并将这些写入（加上自己的写入）转发给后序节点。另一种流行的拓扑是星形结构<sup>注5</sup>：一个指定的根节点将写入转发给所有其他节点。星形拓扑还可以推广到树状结构。

在环形和星形拓扑中，写请求需要通过多个节点才能到达所有的副本，即中间节点需要转发从其他节点收到的数据变更。为防止无限循环，每个节点需要赋予一个唯一的标识符，在复制日志中的每个写请求都标记了已通过的节点标识符<sup>[43]</sup>。如果某个节点收到了包含自身标识符的数据更改，表明该请求已经被处理过，因此会忽略此变更请求，避免重复转发。

环形和星形拓扑的问题是，如果某一个节点发生了故障，在修复之前，会影响其他节点之间复制日志的转发。可以采用重新配置拓扑结构的方法暂时排除掉故障节点。在大多数部署中，这种重新配置必须手动完成。而对于链接更密集的拓扑（如全部到全部），消息可以沿着不同的路径传播，避免了单点故障，因而有更好的容错性。

但另一方面，全链接拓扑也存在一些自身的问题。主要是存在某些网络链路比其他链路更快的情况（例如由于不同网络拥塞），从而导致复制日志之间的覆盖，如图5-9所示。

在图5-9中，客户端A向主节点1的表中首先插入一行，然后客户端B在主节点3上对该行记录进行更新。而在主节点2上，由于网络原因可能出现意外的写日志复制顺序，例如它先接收到了主节点3的更新日志（从主节点2的角度来看，这是对数据库中不存在行的更新操作），之后才接收到主节点1的插入日志（按道理应该在更新日志之前到达）。

注5： 注意不要和前面的星形数据模型混淆，这里主要针对的是节点间的通信模型。

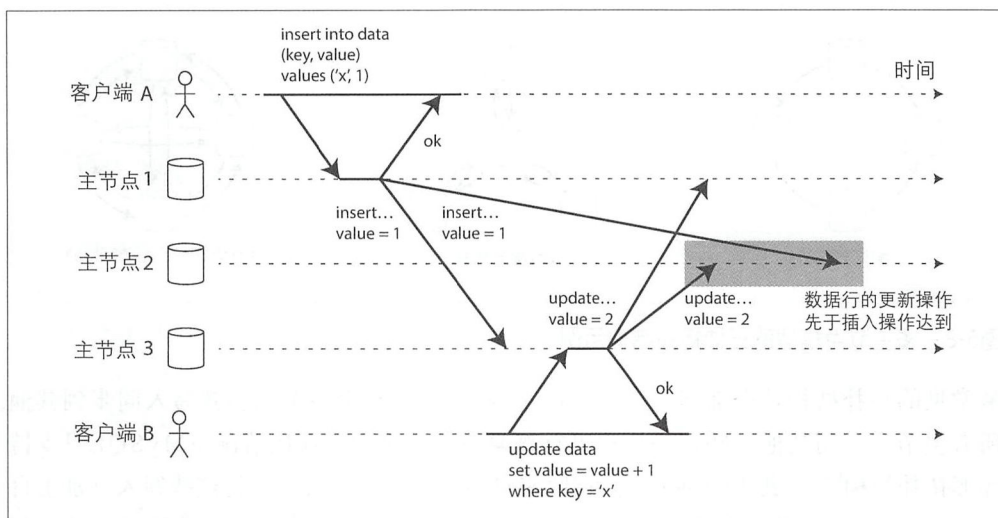


图5-9：对于多主节点复制，某些副本上可能会出现错误的写请求到达顺序

这里涉及到一个因果关系问题，类似于在本章前面“前缀一致读”所看到的：更新操作一定是依赖于先前完成的插入，因此我们要确保所有节点上一定先接收插入日志，然后再处理更新。在每笔写日志里简单地添加时间戳还不够，主要因为无法确保时钟完全同步，因而无法在主节点2上正确地排序所收到日志（参见第8章）。

为了使得日志消息正确有序，可以使用一种称为版本向量的技术，本章稍后将讨论这种技术（参见本章后面的“检测并发写入”）。需要指出，冲突检测技术在许多多主节点复制系统中的实现还不够完善，例如在撰写本书时，PostgreSQL BDR尚不支持写操作的因果排序<sup>[27]</sup>，而MySQL的Tungsten Replicator 甚至还没有基本的冲突检测功能<sup>[34]</sup>。

如果正在使用支持多主节点复制的系统，这些问题都值得注意，仔细查阅相关文档，详细测试这些数据库，以确保它确实提供所期望的功能。

## 无主节点复制

到目前为止本章所讨论的复制方法，包括单主节点和多主节点复制，都是基于这样一种核心思路，即客户端先向某个节点（主节点）发送写请求，然后数据库系统负责将写请求复制到其他副本。由主节点决定写操作的顺序，从节点按照相同的顺序来应用主节点所发送的写日志。

一些数据存储系统则采用了不同的设计思路：选择放弃主节点，允许任何副本直接接

受来自客户端的写请求。其实最早的数据复制系统就是无主节点的（或称为去中心复制，无中心复制）<sup>[1,44]</sup>，但后来到了关系数据库主导的时代，这个想法被大家选择性遗忘了。当亚马逊内部采用了Dynamo系统之后，无主复制又再次成为一种时髦的数据库架构<sup>[37]</sup> 注6。Riak、Cassandra和Voldemort都是受Dynamo启发而设计的无主节点、开源数据库系统，这类数据库也被称为Dynamo风格数据库。

对于某些无主节点系统实现，客户端直接将其写请求发送到多副本，而在其他一些实现中，由一个协调者节点代表客户端进行写入，但与主节点的数据库不同，协调者并不负责写入顺序的维护。我们很快就会看到，这种设计上的差异对数据库的使用方式有着深刻的影响。

## 节点失效时写入数据库

假设一个三副本数据库，其中一个副本当前不可用（例如正在重启以安装系统更新）。在基于主节点复制模型下，如果要继续处理写操作，则需要执行切换操作（参阅本章前面的“处理节点失效”）。

对于无主节点配置，则不存在这样的切换操作。图5-10展示了所发生的情况：用户1234将写请求并行发送到三个副本，有两个可用副本接受写请求，而不可用的副本无法处理该写请求。如果假定三个副本中有两个成功确认写操作，用户1234收到两个确认的回复之后，即可认为写入成功。客户完全可以忽略其中一个副本无法写入的情况。

现在设想一下，失效的节点之后重新上线，而客户端又开始从中读取内容。由于节点失效期间发生的任何写入在该节点上都尚未同步，因此读取可能会得到过期的数据。

为了解决这个问题，当一个客户端从数据库中读取数据时，它不是向一个副本发送请求，而是并行地发送到多个副本。客户端可能会得到不同节点的不同响应，包括某些节点的新值和某些节点的旧值。可以采用版本号技术确定哪个值更新（参见本章后面的“检测并发写入”）。

## 读修复与反熵

复制模型应确保所有数据最终复制到所有的副本。当一个失效的节点重新上线之后，它如何赶上中间错过的那些写请求呢？

---

注6： 亚马逊并不开放使用Dynamo，在其AWS公有云上提供的是DynamoDB，后者采用了完全不同的架构，它是个单主节点的主从复制系统。



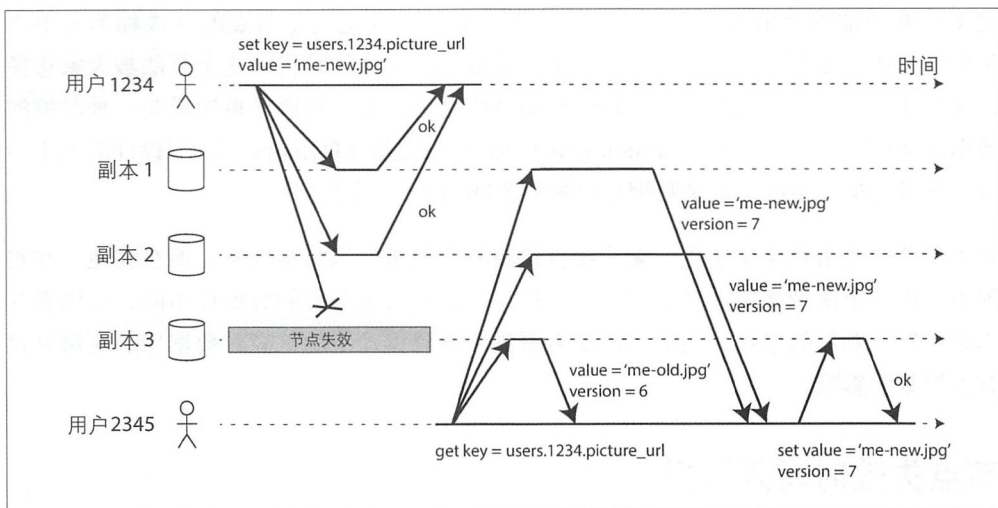


图5-10: quorum读/写, 以及节点失效之后进行读修复

Dynamo风格的数据存储系统经常使用以下两种机制:

### 读修复

当客户端并行读取多个副本时, 可以检测到过期的返回值。例如, 在图5-10中, 用户2345从副本3获得的是版本6, 而从副本1和2得到的是版本7。客户端可以判断副本3一个过期值, 然后将新值写入到该副本。这种方法主要适合那些被频繁读取的场景。

### 反熵过程

此外, 一些数据存储有后台进程不断查找副本之间数据的差异, 将任何缺少的数据从一个副本复制到另一个副本。与基于主节点复制的复制日志不同, 此反熵过程并不保证以特定的顺序复制写入, 并且会引入明显的同步滞后。

并不是所有的系统都实现了上述两种方案。例如, Voldemort目前没有反熵过程。请注意, 当缺少反熵过程的支持时, 由于读时修复只在发生读取时才可能执行修复, 那些很少访问的数据有可能在某些副本中已经丢失而无法检测到, 从而降低了写的持久性。

### 读写quorum

图5-10的例子中, 三个副本中如果有两个以上完成处理, 写入即可认为成功。如果三个副本中只有一个完成了写请求, 会怎样呢? 依次类推, 究竟多少个副本完成才可以认为写成功?

我们知道，成功的写操作要求三个副本中至少两个完成，这意味着至多有一个副本可能包含旧值。因此，在读取时需要至少向两个副本发起读请求，通过版本号可以确定一定至少有一个包含新值。如果第三个副本出现停机或响应缓慢，则读取仍可以继续并返回最新值。

把上述道理推广到一般情况，如果有 $n$ 个副本，写入需要 $w$ 个节点确认，读取必须至少查询 $r$ 个节点，则只要  $w + r > n$ ，读取的节点中一定会包含最新值。例如在前面的例子中， $n = 3$ ， $w = 2$ ， $r = 2$ 。满足上述这些 $r$ 、 $w$ 值的读/写操作称之为法定票数读（或仲裁读）或法定票数写（或仲裁写）<sup>[44]注7</sup>。也可以认为 $r$ 和 $w$ 是用于判定读、写是否有效的最低票数。

在Dynamo风格的数据库中，参数 $n$ 、 $w$ 和 $r$ 通常是可配置的。一个常见的选择是设置 $n$ 为某奇数（通常为3或5）， $w = r = (n + 1) / 2$ （向上舍入）。也可以根据自己的需求灵活调整这些配置。例如，对于读多写少的负载，设置 $w = n$ 和 $r = 1$ 比较合适，这样读取速度更快，但是一个失效的节点就会使得数据库所有写入因无法完成quorum而失败。



集群中可能存在多于 $n$ 个节点，但是数据只会保存在所设定的 $n$ 个节点上。我们可以对数据集进行分区，从而支持比节点容纳上限更大的数据集，第6章我们将讨论分区技术。

仲裁条件 $w + r > n$ 定义了系统可容忍的失效节点数，如下所示：

- 当 $w < n$ ，如果一个节点不可用，仍然可以处理写入。
- 当 $r < n$ ，如果一个节点不可用，仍然可以处理读取。
- 假定 $n = 3$ ， $w = 2$ ， $r = 2$ ，则可以容忍一个不可用的节点。
- 假定 $n = 5$ ， $w = 3$ ， $r = 3$ ，则可以容忍两个不可用的节点。如图5-11所示。
- 通常，读取和写入操作总是并行发送到所有的 $n$ 个副本。参数 $w$ 和参数 $r$ 只是决定要等待的节点数。即有多少个节点需要返回结果，我们才能判断出结果的正确性。

如果可用节点数小于所需的 $w$ 或 $r$ ，则写入或读取就会返回错误。不可用的原因可能有很多种，包括节点崩溃或者断电而关机，执行操作时出错（例如磁盘已满而无法写入），客户端和节点之间的网络中断等。这里，我们只需关心节点是否有返回值，而不需区分出错的具体原因。

---

注7： 有时也被称为严格法定票数，与之对应的称为宽松法定票数，参见本章后面的“宽松的quorum与数据回传”。

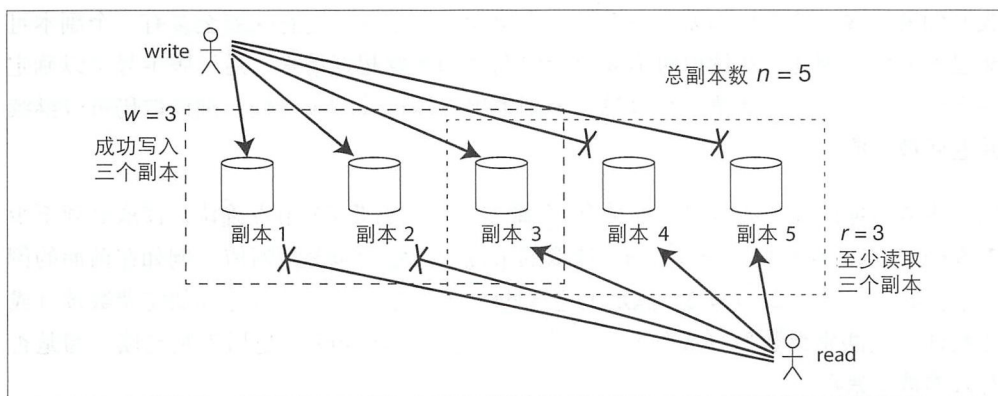


图5-11：如果  $w + r > n$ ，则并发  $r$  个读请求中一定包含最新值

## Quorum一致性的局限性

如果有  $n$  个副本，并且配置  $w$  和  $r$ ，使得  $w + r > n$ ，可以预期可以读取到一个最新值。之所以这样，是因为成功写入的节点集合和读取的节点集合必然有重合，这样读取的节点中至少有一个具有最新值（见图5-11）。

通常，设定  $r$  和  $w$  为简单多数（多于  $n / 2$ ）节点，即可确保  $w + r > n$ ，且同时容忍多达  $n / 2$  个节点故障。但是，quorum 不一定非得是多数，读和写的节点集中有一个重叠的节点才是最关键的。设定其他的 quorum 分配数也是可行的<sup>[45]</sup>。

也可以将  $w$  和  $r$  设置为较小的数字，从而  $w + r \leq n$ （即不满足仲裁条件）。此时，读取和写入操作仍会被发送到  $n$  个节点，但只需等待更少的节点回应即可返回。

由于  $w$  和  $r$  配置的节点数较小，读取请求当中可能恰好没有包含新值的节点，因此最终可能会返回一个过期的旧值。好的一方面是，这种配置可以获得更低的延迟和更高的可用性，例如网络中断，许多副本变得无法访问，相比而言有更高的概率继续处理读取和写入。只有当可用的副本数已经低于  $w$  或  $r$  时，数据库才会变得无法读/写，即处于不可用状态。

即使在  $w + r > n$  的情况下，也可能存在返回旧值的边界条件。这主要取决于具体实现，可能的情况包括：

- 如果采用了 sloppy quorum（参阅本章后面的“宽松的 quorum 与数据回传”），写操作的  $w$  节点和读取的  $r$  节点可能完全不同，因此无法保证读写请求一定存在重叠的节点<sup>[46]</sup>。



- 如果两个写操作同时发生，则无法明确先后顺序。这种情况下，唯一安全的解决方案是合并并发写入（参见本章前面的“处理写冲突”）。如果根据时间戳（最后写入获胜）挑选胜者，则由于时钟偏差问题<sup>[35]</sup>，某些写入可能会被错误地抛弃。
- 如果写操作与读操作同时发生，写操作可能仅在一部分副本上完成。此时，读取时返回旧值还是新值存在不确定性。
- 如果某些副本上已经写入成功，而其他一些副本发生写入失败（例如磁盘已满），且总的成功副本数少于 $w$ ，那些已成功副本上不会做回滚。这意味着尽管这样的写操作被视为失败，后续的读操作仍可能返回新值<sup>[47]</sup>。
- 如果具有新值的节点后来发生失效，但恢复数据来自某个旧值，则总的新值副本数会低于 $w$ ，这就打破了之前的判定条件。
- 即使一切工作正常，也会出现一些边界情况，如第9章所介绍的“可线性化与quorum”。

因此，虽然quorum设计上似乎可以保证读取最新值，但现实情况却往往更加复杂。Dynamo风格的数据库通常是针对最终一致性场景而优化的。我们建议最好不要把参数 $w$ 和 $r$ 视为绝对的保证，而是一种灵活可调的读取新值的概率。

例如，这里通常无法得到本章前面的“复制滞后问题”中所罗列的一致性保证，包括写后读、单调读、前缀一致读等，因此前面讨论种种异常同样会发生在这里。如果确实需要更强的保证，需要考虑事务与共识问题，接下来的第7章和第9章将对此展开讨论。

## 监控旧值

从运维角度来看，监视数据库是否返回最新结果非常重要。即使应用程序可以容忍读取旧值，也需要仔细了解复制的当前运行状态。如果已经出现了明显的滞后，它就是个重要的信号提醒我们需要采取必要措施来排查原因（例如网络问题或节点超负荷）。

对于主从复制的系统，数据库通常会导出复制滞后的相关指标，可以将其集成到统一监控模块。原理大概是这样，由于主节点和从节点上写入都遵从相同的顺序，而每个节点都维护了复制日志执行的当前偏移量。通过对比主节点和从节点当前偏移量的差值，即可衡量该从节点落后于主节点的程度。

然而，对于无主节点复制的系统，并没有固定的写入顺序，因而监控就变得更加困

难。而且，如果数据库只支持读时修复（不支持反熵），那么旧值的落后就没有一个上限。例如如果一个值很少被访问，那么所返回的旧值可能非常之古老。

目前针对无主节点复制系统已经有一些研究，根据参数 $n$ ， $w$ 和 $r$ 来预测读到旧值的期望百分比<sup>[48]</sup>。不过，总体讲还不是很普及。即便如此，将旧值监控纳入到数据库标准指标集中还是很有必要。要知道，最终一致性其实是个非常模糊的保证，从可操作性上讲，量化究竟何为“最终”很有实际价值。

## 宽松的quorum与数据回传

配置适当quorum的数据库系统可以容忍某些节点故障，也不需要执行故障切换。它们还可以容忍某些节点变慢，这是因为请求并不需要等待所有 $n$ 个节点的响应，只需 $w$ 或 $r$ 节点响应即可。对于需要高可用和低延迟的场景来说，还可以容忍偶尔读取旧值，所有这些特性使之具有很高的吸引力。

但是，quorum并不总如期待的那样提供高容错能力。一个网络中断可以很容易切断一个客户端到多数数据库节点的链接。尽管这些集群节点是活着的，而且其他客户端也确实可以正常链接，但是对于断掉链接的客户端来讲，情况无疑等价于集群整体失效。这种情况下，很可能无法满足最低的 $w$ 和 $r$ 所要求的节点数，因此导致客户端无法满足quorum要求。

在一个大规模集群中（节点数远大于 $n$ 个），客户可能在网络中断期间还能连接到某些数据库节点，但这些节点又不是能够满足数据仲裁的那些节点。此时，数据库设计者就面临着一个选择：

- 如果无法达到 $w$ 或 $r$ 所要求quorum，将错误明确地返回给客户端？
- 或者，我们是否应该接受该写请求，只是将它们暂时写入一些可访问的节点中？注意，这些节点并不在 $n$ 个节点集合中。

后一种方案称之为放松的仲裁<sup>[37]</sup>：写入和读取仍然需要 $w$ 和 $r$ 个成功的响应，但包含了那些并不在先前指定的 $n$ 个节点。打个比方，如果你把不小心把自己锁在房子外面，可能会敲开邻居家的门，请求否可以坐在沙发上暂时休息一下。

一旦网络问题得到解决，临时节点需要把接收到的写入全部发送到原始主节点上。这就是所谓的数据回传（或暗示移交）。即一旦你找到了房子的钥匙，你的邻居会礼貌地请你离开沙发回到自己的家中。

可以看出，sloppy quorum对于提高写入可用性特别有用：只要有任何 $w$ 个节点可用，

数据库就可以接受新的写入。然而这意味着，即使满足  $w + r > n$ ，也不能保证在读取某个键时，一定能读到最新值，因为新值可能被临时写入  $n$  之外的某些节点且尚未回传过来<sup>[47]</sup>。

因此，sloppy quorum并非传统意义上quorum。而更像是为了数据持久性而设计的一个保证措施，除非回传结束，否则它无法保证客户端一定能从  $r$  个节点读到新值。

目前，所有Dynamo风格的系统都已经支持sloppy quorum。在Riak中，默认启用，而在Cassandra和Voldemort中则默认关闭<sup>[46,49,50]</sup>。

## 多数据中心操作

我们之前以多数据中心为例介绍了多主节点复制（参见本章前面的“多主节点复制”）。而无主节点复制由于旨在更好地容忍并发写入冲突，网络中断和延迟尖峰等，因此也可适用于多数据中心操作。

Cassandra和Voldemort在其默认配置的无主节点模型中都支持跨数据中心操作：副本的数量  $n$  是包含所有数据中心的节点总数。配置时，可以指定每个数据中心各有多少副本。每个客户端的写入都会发送到所有副本，但客户端通常只会等待来自本地数据中心内的quorum节点数的确认，这样避免了高延迟和跨数据中心可能的网络异常。尽管可以灵活配置，但对远程数据中心的写入由于延迟很高，通常都被配置为异步方式<sup>[50,51]</sup>。

Riak则将客户端与数据库节点之间的通信限制在一个数据中心内，因此  $n$  描述的是一个数据中心内的副本数量。集群之间跨数据中心的复制则在后台异步运行，类似于多主节点复制风格<sup>[52]</sup>。

## 检测并发写

Dynamo风格的数据库允许多个客户端对相同的主键同时发起写操作，即使采用严格的quorum机制也可能会发生写冲突。这与多主节点复制类似（参见本章前面的“处理写冲突”），此外，由于读时修复或者数据回传也会导致并发写冲突。

一个核心问题是，由于网络延迟不稳定或者局部失效，请求在不同的节点上可能会呈现不同的顺序。如图5-12所示，对于包含三个节点的数据系统，客户端A和B同时向主键X发起写请求：

- 节点1收到来自客户端A的写请求，但由于节点失效，没有收到客户端B的写请求。



- 节点2首先收到A的写请求，然后是B的写请求。
- 与节点2相反，节点3首先收到B的写请求，然后是A的写请求。

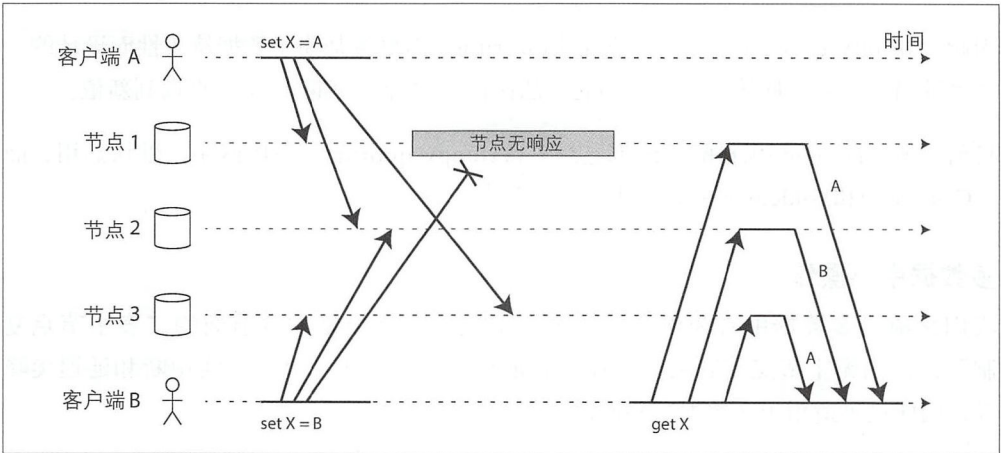


图5-12: Dynamo风格的数据库在并发写入时缺乏顺序保证

如果节点每当收到新的写请求时就简单地覆盖原有的主键，那么这些节点将永久无法达成一致，如图5-12中的所示，节点2认为X的最终值是B，而其他节点认为值是A。

我们知道副本应该收敛于相同的内容，这样才能达成最终一致。但如何才能做到呢？有人可能希望数据副本之间能自动处理，然而非常不幸，目前大多数的系统实现都无法令人满意，如果你不想丢失数据，应用开发者必须了解很多关于数据库内部冲突处理的机制。

我们已经在本章前面的“处理写冲突”简要介绍了一些解决冲突的技巧。在总结本章之前，我们来更详细地探讨这个问题。

### 最后写入者获胜（丢弃并发写入）

一种实现最终收敛的方法是，每个副本总是保存最新值，允许覆盖并丢弃旧值。那么，假定每个写请求都最终同步到所有副本，只要我们有一个明确的方法来确定哪一个写入是最新的，则副本可以最终收敛到相同的值。

这个想法其实有些争议，关键点在于前面所提到关于如何定义“最新”。在图5-12的例子中，当客户端向数据库节点发送写请求时，一个客户端无法意识到另一个客户端，也不清楚哪一个先发生。其实，争辩哪个先发生没有太大意义，当我们说支持写入并发，也就意味着它们的顺序是不确定的。

即使无法确定写请求的“自然顺序”，我们可以强制对其排序。例如，为每个写请求附加一个时间戳，然后选择最新即最大的时间戳，丢弃较早时间戳的写入。这个冲突解决算法被称为最后写入者获胜（last write wins, LWW），它是Cassandra<sup>[53]</sup>仅有的冲突解决方法，而在Riak<sup>[35]</sup>中，它是可选方案之一。

LWW可以实现最终收敛的目标，但是以牺牲数据持久性为代价。如果同一个主键有多个并发写，即使这些并发写都向客户端报告成功（因为完成了写入 $w$ 个副本），但最后只有一个写入值会存活下来，其他的将被系统默默丢弃。此外，LWW甚至可能会删除那些非并发写，我们将在第8章“时间戳与事件顺序”中举例说明。

在一些场景如缓存系统，覆盖写是可以接受的。如果覆盖、丢失数据不可接受，则LWW并不是解决冲突很好的选择。

要确保LWW安全无副作用的唯一方法是，只写入一次然后写入值视为不可变，这样就避免了对同一个主键的并发（覆盖）写。例如，Cassandra的一个推荐使用方法就是采用UUID作为主键，这样每个写操作都针对的不同的、系统唯一的主键<sup>[53]</sup>。

## Happens-before关系和并发

如何判断两个操作是否是并发呢？首先为了建立起一个快速的直觉判断，我们先来看一些例子：

- 图5-9中，两个写入不是并发的：A的插入操作发生在B的增量修改之前，B的递增是基于A插入的行。换句话说，B后发生，其操作建立在A基础之上。A和B属于因果依赖关系。
- 另一个例子，图5-12中的两个写入则是并发的：每个客户端启动写操作时，并不知道另一个客户端是否也在同一个主键上执行操作。因此，操作之间不存在因果关系。

如果B知道A，或者依赖于A，或者以某种方式在A基础上构建，则称操作A在操作B之前发生。这是定义何为并发的关键。事实上，我们也可以简单地说，如果两个操作都不在另一个之前发生，那么操作是并发的（或者两者都不知道对方）<sup>[54]</sup>。

因此，对于两个操作A和B，一共存在三种可能性：A在B之前发生，或者B在A之前发生，或者A和B并发。我们需要的是一个算法来判定两个操作是否并发。如果一个操作发生在另一个操作之前，则后面的操作可以覆盖较早的操作。如果属于并发，就需要解决潜在的冲突问题。

## 并发性、时间和相对性

通常如果两个操作“同时”发生，则称之为并发，然而事实上，操作是否在时间上重叠并不重要。由于分布式系统中复杂的时钟同步问题（第8章将会详细讨论），现实当中，我们很难严格确定它们是否同时发生。

为更好地定义并发性，我们并不依赖确切的发生时间，即不管物理的时机如何，如果两个操作并不需要意识到对方，我们即可声称它们是并发操作。一些人尝试把这个思路与物理学中狭义相对论联系起来<sup>[54]</sup>，后者引入了“信息传递不能超越光速”的假定，如果两个事件发生的间隔短于光在它们之间的折返，那么这两个事件不可能有相互影响，因此就是并发。

在计算机系统中，即使光速快到允许一个操作影响到另一个操作，但两个操作仍可能被定义为并发。例如，发生了网络拥塞或中断，可能就会出现两个操作由于网络问题导致一个操作无法感知另一个，因此二者成为并发。

### 确定前后关系

我们来看一个确定操作并发性的算法，即两个操作究竟属于并发还是一个发生在另一个之前（依赖关系）。简单起见，我们先从只有一个副本的数据库开始，在阐明其原理之后，将其推广到有多个副本的无主节点数据库。

图5-13的例子是两个客户端同时向购物篮车加商品（如果觉得这个例子太微不足道，可以类比为，两个空中交警员同时把飞机添加到他们所管理的追踪目标里）。初始时购物车为空。然后两个客户端向数据库共发出五次写入操作：

1. 客户端1首先将牛奶加入购物车。这是第一次写入该主键的值，服务器保存成功后分配版本1，服务器将值与版本号一起返回给该客户端1。
2. 客户端2将鸡蛋加入购物车，此时它并不知道客户端1已添加了牛奶，而是认为鸡蛋是购物车中的唯一物品。服务器为此写入并分配版本2，然后将鸡蛋和牛奶存储为两个单独的值，最好将这两个值与版本号2返回给客户端2。
3. 同理，客户端1也并不意味着上述步骤2，想要将面粉加入购物车，且以为购物车的内容应该是[牛奶，面粉]，将此值与版本号1一起发送到服务器。服务器可以从版本号中知道[牛奶，面粉]的新值要取代先前值[牛奶]，但值[鸡蛋]则是新的并发操作。因此，服务器将版本3分配给[牛奶，面粉]并覆盖版本1的[牛奶]，同时保留版本2的值[鸡蛋]，将二者同时返回给客户端1。



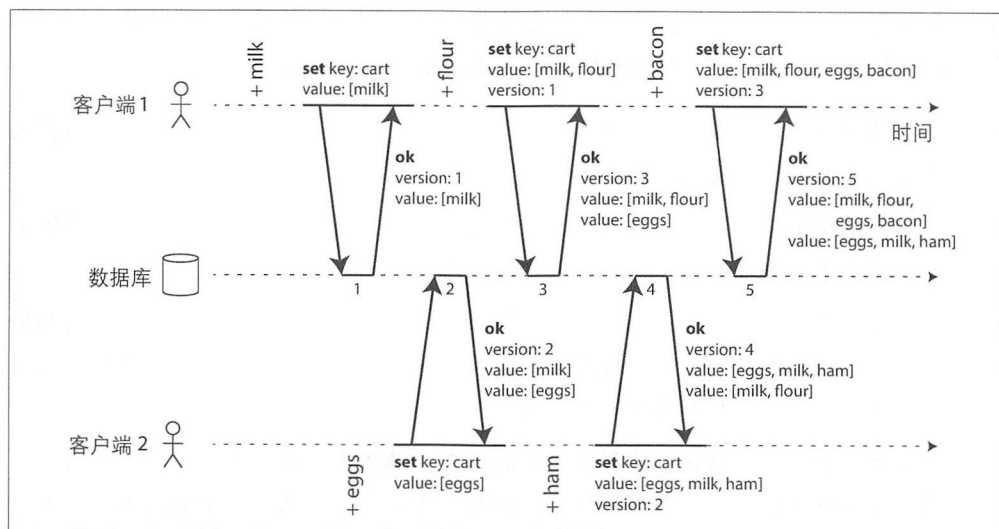


图5-13：以并发操作购物车为例说明捕获事件的因果关系

- 同时，客户端2想要加入火腿，也不知道客户端1刚刚加了面粉。其在最后一个响应中从服务器收到的两个值是[牛奶]和[蛋]，现在合并这些值，并添加火腿形成一个新的值[鸡蛋，牛奶，火腿]。它将该值与前一个版本号2一起发送到服务器。服务器检测到版本2会覆盖[鸡蛋]，但与[牛奶，面粉]是同时发生，所以设置为版本4并将所有这些值发送给客户端2。
- 最后，客户端1想要加培根。它以前在版本3中从服务器接收[牛奶，面粉]和[鸡蛋]，所以合并这些值，添加培根，并将最终值[牛奶，面粉，鸡蛋，培根]连同版本号3来覆盖[牛奶，面粉]，但与[鸡蛋，牛奶，火腿]并发，所以服务器会保留这些并发值。

图5-13操作之间的数据流可以通过图5-14形象展示。箭头表示某个操作发生在另一个操作之前，即后面的操作“知道”或是“依赖”于前面的操作。在这个例子中，因为总有另一个操作同时进行，所以每个客户端都没有时时刻刻和服务器上的数据保持同步。但是，新版本值最终会覆盖旧值，且不会发生已写入值的丢失。

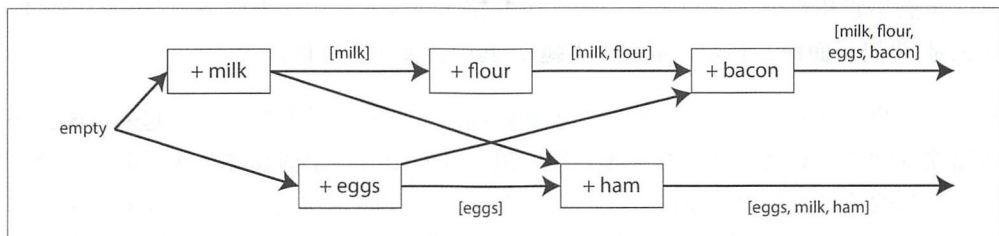


图5-14：上图5-13背后的因果关系

需要注意的是，服务器判断操作是否并发的依据主要依靠对比版本号，而并不需要解释新旧值本身（值可以是任何数据结构）。算法的工作流程如下：

- 服务器为每个主键维护一个版本号，每当主键新值写入时递增版本号，并将新版本号与写入的值一起保存。
- 当客户端读取主键时，服务器将返回所有（未被覆盖的）当前值以及最新的版本号。且要求写之前，客户必须先发送读请求。
- 客户端写主键，写请求必须包含之前读到的版本号、读到的值和新值合并后的集合。写请求的响应可以像读操作一样，会返回所有当前值，这样就可以像购物车例子那样一步步链接起多个写入的值。
- 当服务器收到带有特定版本号的写入时，覆盖该版本号或更低版本的所有值（因为知道这些值已经被合并到新传入的值集合中），但必须保存更高版本号的所有值（因为这些值与当前的写操作属于并发）。

当写请求包含了前一次读取的版本号时，意味着修改的是基于以前的状态。如果一个写请求没有包含版本号，它将与所有其他写入同时进行，不会覆盖任何已有值，其传入的值将包含在后续读请求的返回值列表当中。

## 合并同时写入的值

上述算法可以保证不会发生数据丢弃，但不幸的是，客户端需要做一些额外的工作：即如果多个操作并发发生，则客户端必须通过合并并发写入的值来继承旧值。Riak称这些并发值为兄弟关系。

合并本质上与先前讨论的多节点复制冲突解决类似（参阅本章前面的“处理写冲突”）。一个简单的方法是基于版本号或时间戳（即最后写入获胜）来选择其中的一个值，但这意味着会丢失部分数据。所以，需要在应用程序代码中额外做些工作。

以购物车为例，合并并发值的合理方式是包含新值和旧值（union操作）。图5-14中，两个客户端最后的值分别是 [牛奶，面粉，鸡蛋，熏肉]和[鸡蛋，牛奶，火腿]。注意，虽然牛奶和鸡蛋只是写入了一次，但它在两个客户端中均有出现。合并的最终值应该是[牛奶，面粉，鸡蛋，培根，火腿]，其中去掉了重复值。

然而，设想一下人们也可以在购物车中删除商品，此时把并发值都合并起来可能会导致错误的结果：如果合并了两个客户端的值，且其中有一个商品被某客户端删除掉，则被删除的项目会再次出现在合并的终值中<sup>[37]</sup>。为了防止该问题，项目在删除时不能简单地从数据库中删除，系统必须保留一个对应的版本号以恰当的标记该项目需要在

合并时被剔除。这种删除标记被称为墓碑（之前我们在第3章“哈希索引”日志压缩时提到过）。

考虑到在应用代码中合并非常复杂且容易出错，因此可以设计一些专门的数据结构来自动执行合并，例如，Riak支持称为CRDT一系列数据结构<sup>[38,39,55]</sup>（具体参见本章前面的“自动冲突解决”），以合理的方式高效自动合并，包括支持删除标记。

## 版本矢量

图5-13中的示例只有一个副本。如果存在多个副本但没有主节点，算法又该如何呢？

图5-13使用单个版本号来捕获操作之间的依赖关系，当多个副本同时接受写入时，这是不够的。因此我们需要为每个副本和每个主键均定义一个版本号。每个副本在处理写入时增加自己的版本号，并且跟踪从其他副本看到的版本号。通过这些信息来指示要覆盖哪些值、该保留哪些并发值。

所有副本的版本号集合称为版本矢量<sup>[56]</sup>。这种思路还有一些变体，但最有趣的可能是在Riak 2.0<sup>[58,59]</sup>中使用的虚线版本矢量<sup>[57]</sup>。我们无法在此深入其细节，但是它的工作方式与购物车例子所展示的非常相似。

与图5-13中版本号类似，当读取数据时，数据库副本会返回版本矢量给客户端，而在随后写入时需要将版本信息包含在请求当中一起发送到数据库。Riak将版本矢量编码为一个称之为因果上下文的字符串。版本矢量技术使数据库可以区分究竟应该覆盖写还是保留并发值。

另外，就像单副本的例子一样，应用程序仍然需要执行合并操作。版本矢量可以保证从某一个副本读取值然后写入到另一个副本，而这些值可能会导致在其他副本上衍生出来新的“兄弟”值，但至少不会发生数据丢失且可以正确合并所有并发值。



### 版本矢量和矢量时钟

版本矢量有时也被称为矢量时钟，然而两者并不完全相同，细微差别可参阅文献<sup>[57,60,61]</sup>。简而言之，当需要比较副本状态时，应当采用版本矢量。

## 小结

本章，我们详细探讨了复制相关的话题。复制或者多副本技术主要服务于以下目的：



- 高可用性：即使某台机器（或多台机器，或整个数据中心）出现故障，系统也能保持正常运行。
- 连接断开与容错：允许应用程序在出现网络中断时继续工作。
- 低延迟：将数据放置在距离用户较近的地方，从而实现更快地交互。
- 可扩展性：采用多副本读取，大幅提高系统读操作的吞吐量。

在多台机器上保存多份相同的数据副本，看似只是个很简单的目标，但事实上复制技术是一个非常烧脑的问题。需要仔细考虑并发以及所有可能出错的环节，并小心处理故障之后的各种情形。最最基本的，要处理好节点不可用与网络中断问题，这里甚至还没考虑一些更隐蔽的失效场景，例如由于软件bug而导致的无提示的数据损坏。

我们主要讨论了三种多副本方案：

### 主从复制

所有的客户端写入操作都发送到某一个节点（主节点），由该节点负责将数据更改事件发送到其他副本（从节点）。每个副本都可以接收读请求，但内容可能是过期值。

### 多主节点复制

系统存在多个主节点，每个都可以接收写请求，客户端将写请求发送到其中的一个主节点上，由该主节点负责将数据更改事件同步到其他主节点和自己的从节点。

### 无主节点复制

客户端将写请求发送到多个节点上，读取时从多个节点上并行读取，以此检测和纠正某些过期数据。

每种方法都有其优点和缺点。主从复制非常流行，主要是因为它很容易理解，也不需要担心冲突问题。而万一出现节点失效、网络中断或者延迟抖动等情况，多主节点和无主节点复制方案会更加可靠，不过背后的代价则是系统的复杂性和弱一致性保证。

复制可以是同步的，也可以是异步的，而一旦发生故障，二者的表现差异会对系统行为产生深远的影响。在系统稳定状态下异步复制性能优秀，但仍须认真考虑一旦出现复制滞后和节点失效两种场景会导致何种影响。万一某个主节点发生故障，而一个异步更新的从节点被提升为新的主节点，要意识到最新确认的数据可能有丢失的风险。

我们还分析了由于复制滞后所引起的一些奇怪效应，并讨论了以下一致性模型，来帮助应用程序处理复制滞后：

### 写后读一致性

保证用户总能看到自己所提交的最新数据。

### 单调读

用户在某个时间点读到数据之后，保证此后不会出现比该时间点更早的数据。

### 前缀一致读

保证数据之间的因果关系，例如，总是以正确的顺序先读取问题，然后看到回答。

最后，我们讨论了多主节点和无主节点复制方案所引入的并发问题。即由于多个写可能同时发生，继而可能产生冲突。为此，我们研究了一个算法使得数据库系统可以判定某操作是否发生在另一个操作之前，或者是同时发生。接下来，探讨采用合并并发更新值的方法来解决冲突。

下一章我们继续研究多节点上数据的分布问题，与本章不同的是，它是针对一个大型数据集而采用分区技术。

## 参考文献

- [1] Bruce G. Lindsay, Patricia Griffiths Selinger, C. Galtieri, et al.: “Notes on Distributed Databases,” IBM Research, Research Report RJ2571(33471), July 1979.
- [2] “Oracle Active Data Guard Real-Time Data Protection and Availability,” Oracle White Paper, June 2013.
- [3] “AlwaysOn Availability Groups,” in *SQL Server Books Online*, Microsoft, 2012.
- [4] Lin Qiao, Kapil Surlaker, Shirshanka Das, et al.: “On Brewing Fresh Espresso: LinkedIn’s Distributed Data Serving Platform,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2013.
- [5] Jun Rao: “Intra-Cluster Replication for Apache Kafka,” at *ApacheCon North America*, February 2013.
- [6] “Highly Available Queues,” in *RabbitMQ Server Documentation*, Pivotal Software, Inc., 2014.
- [7] Yoshinori Matsunobu: “Semi-Synchronous Replication at Facebook,” *yoshinorimatsunobu.blogspot.co.uk*, April 1, 2014.

- [8] Robbert van Renesse and Fred B. Schneider: “Chain Replication for Supporting High Throughput and Availability,” at *6th USENIX Symposium on Operating System Design and Implementation (OSDI)*, December 2004.
- [9] Jeff Terrace and Michael J. Freedman: “Object Storage on CRAQ: High- Throughput Chain Replication for Read-Mostly Workloads,” at *USENIX Annual Technical Conference (ATC)*, June 2009.
- [10] Brad Calder, Ju Wang, Aaron Ogus, et al.: “Windows Azure Storage: A Highly Available Cloud Storage Service with Strong Consistency,” at *23rd ACM Symposium on Operating Systems Principles (SOSP)*, October 2011.
- [11] Andrew Wang: “Windows Azure Storage,” *umbrant.com*, February 4, 2016.
- [12] “Percona Xtrabackup - Documentation,” Percona LLC, 2014.
- [13] Jesse Newland: “GitHub Availability This Week,” *github.com*, September 14, 2012.
- [14] Mark Imbriaco: “Downtime Last Saturday,” *github.com*, December 26, 2012.
- [15] John Hugg: “ ‘All in’ with Determinism for Performance and Testing in Distributed Systems,” at *Strange Loop*, September 2015.
- [16] Amit Kapila: “WAL Internals of PostgreSQL,” at *PostgreSQL Conference (PGCon)*, May 2012.
- [17] *MySQL Internals Manual*. Oracle, 2014.
- [18] Yogeshwer Sharma, Philippe Ajoux, Petchean Ang, et al.: “Wormhole: Reliable Pub-Sub to Support Geo-Replicated Internet Services,” at *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, May 2015.
- [19] “Oracle GoldenGate 12c: Real-Time Access to Real-Time Information,” Oracle White Paper, October 2013.
- [20] Shirshanka Das, Chavdar Botev, Kapil Surlaker, et al.: “All Aboard the Databus!,” at *ACM Symposium on Cloud Computing (SoCC)*, October 2012.
- [21] Greg Sabino Mullane: “Version 5 of Bucardo Database Replication System,” *blog.endpoint.com*, June 23, 2014.



- [22] Werner Vogels: “Eventually Consistent,” *ACM Queue*, volume 6, number 6, pages 14–19, October 2008. doi:10.1145/1466443.1466448.
- [23] Douglas B. Terry: “Replicated Data Consistency Explained Through Baseball,” Microsoft Research, Technical Report MSR-TR-2011-137, October 2011.
- [24] Douglas B. Terry, Alan J. Demers, Karin Petersen, et al.: “Session Guarantees for Weakly Consistent Replicated Data,” at *3rd International Conference on Parallel and Distributed Information Systems* (PDIS), September 1994. doi:10.1109/PDIS.1994.331722.
- [25] Terry Pratchett: *Reaper Man: A Discworld Novel*. Victor Gollancz, 1991. ISBN: 978-0-575-04979-6.
- [26] “Tungsten Replicator,” Continuent, Inc., 2014.
- [27] “BDR 0.10.0 Documentation,” The PostgreSQL Global Development Group, *bdr-project.org*, 2015.
- [28] Robert Hodges: “If You \*Must\* Deploy Multi-Master Replication, Read This First,” *scale-out-blog.blogspot.co.uk*, March 30, 2012.
- [29] J. Chris Anderson, Jan Lehnardt, and Noah Slater: *CouchDB: The Definitive Guide*. O’Reilly Media, 2010. ISBN: 978-0-596-15589-6.
- [30] AppJet, Inc.: “Etherpad and EasySync Technical Manual,” *github.com*, March 26, 2011.
- [31] John Day-Richter: “What’s Different About the New Google Docs: Making Collaboration Fast,” *googledrive.blogspot.com*, 23 September 2010.
- [32] Martin Kleppmann and Alastair R. Beresford: “A Conflict-Free Replicated JSON Datatype,” arXiv:1608.03960, August 13, 2016.
- [33] Frazer Clement: “Eventual Consistency – Detecting Conflicts,” *messagepassing.blogspot.co.uk*, October 20, 2011.
- [34] Robert Hodges: “State of the Art for MySQL Multi-Master Replication,” at *Percona Live: MySQL Conference & Expo*, April 2013.
- [35] John Daily: “Clocks Are Bad, or, Welcome to the Wonderful World of Distributed

Systems,” *basho.com*, November 12, 2013.

[36] Riley Berton: “Is Bi-Directional Replication (BDR) in Postgres Transactional?,” *sdf.org*, January 4, 2016.

[37] Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, et al.: “Dynamo: Amazon’s Highly Available Key-Value Store,” at *21st ACM Symposium on Operating Systems Principles (SOSP)*, October 2007.

[38] Marc Shapiro, Nuno Preguiça, Carlos Baquero, and Marek Zawirski: “A Comprehensive Study of Convergent and Commutative Replicated Data Types,” INRIA Research Report no. 7506, January 2011.

[39] Sam Elliott: “CRDTs: An UPDATE (or Maybe Just a PUT),” at *RICON West*, October 2013.

[40] Russell Brown: “A Bluffers Guide to CRDTs in Riak,” *gist.github.com*, October 28, 2013.

[41] Benjamin Farinier, Thomas Gazagnaire, and Anil Madhavapeddy: “Mergeable Persistent Data Structures,” at *26es Journées Francophones des Langages Applicatifs (JFLA)*, January 2015.

[42] Chengzheng Sun and Clarence Ellis: “Operational Transformation in Real-Time Group Editors: Issues, Algorithms, and Achievements,” at *ACM Conference on Computer Supported Cooperative Work (CSCW)*, November 1998.

[43] Lars Hofhansl: “HBASE-7709: Infinite Loop Possible in Master/Master Replication,” *issues.apache.org*, January 29, 2013.

[44] David K. Gifford: “Weighted Voting for Replicated Data,” at *7th ACM Symposium on Operating Systems Principles (SOSP)*, December 1979. doi: 10.1145/800215.806583.

[45] Heidi Howard, Dahlia Malkhi, and Alexander Spiegelman: “Flexible Paxos: Quorum Intersection Revisited,” *arXiv:1608.06696*, August 24, 2016.

[46] Joseph Blomstedt: “Re: Absolute Consistency,” email to *riak-users* mailing list, *lists.basho.com*, January 11, 2012.

[47] Joseph Blomstedt: “Bringing Consistency to Riak,” at *RICON West*, October 2012.

- [48] Peter Bailis, Shivaram Venkataraman, Michael J. Franklin, et al.: “Quantifying Eventual Consistency with PBS,” *Communications of the ACM*, volume 57, number 8, pages 93–102, August 2014. doi:10.1145/2632792.
- [49] Jonathan Ellis: “Modern Hinted Handoff,” *datastax.com*, December 11, 2012.
- [50] “Project Voldemort Wiki,” *github.com*, 2013.
- [51] “Apache Cassandra 2.0 Documentation,” DataStax, Inc., 2014.
- [52] “Riak Enterprise: Multi-Datacenter Replication.” Technical whitepaper, Basho Technologies, Inc., September 2014.
- [53] Jonathan Ellis: “Why Cassandra Doesn’t Need Vector Clocks,” *datastax.com*, September 2, 2013.
- [54] Leslie Lamport: “Time, Clocks, and the Ordering of Events in a Distributed System,” *Communications of the ACM*, volume 21, number 7, pages 558–565, July 1978. doi:10.1145/359545.359563.
- [55] Joel Jacobson: “Riak 2.0: Data Types,” *blog.joeljacobson.com*, March 23, 2014.
- [56] D. Stott Parker Jr., Gerald J. Popek, Gerard Rudisin, et al.: “Detection of Mutual Inconsistency in Distributed Systems,” *IEEE Transactions on Software Engineering*, volume 9, number 3, pages 240–247, May 1983. doi:10.1109/TSE.1983.236733.
- [57] Nuno Preguiça, Carlos Baquero, Paulo Sérgio Almeida, et al.: “Dotted Version Vectors: Logical Clocks for Optimistic Replication,” arXiv:1011.5808, November 26, 2010.
- [58] Sean Cribbs: “A Brief History of Time in Riak,” at *RICON*, October 2014.
- [59] Russell Brown: “Vector Clocks Revisited Part 2: Dotted Version Vectors,” *basho.com*, November 10, 2015.
- [60] Carlos Baquero: “Version Vectors Are Not Vector Clocks,” *haslab.wordpress.com*, July 8, 2011.
- [61] Reinhard Schwarz and Friedemann Mattern: “Detecting Causal Relationships in Distributed Computations: In Search of the Holy Grail,” *Distributed Computing*, volume 7, number 3, pages 149–174, March 1994. doi:10.1007/BF02277859.





# 数据分区

我们必须摆脱串行的限制。明确状态定义，提供优先级与属性描述。我们必须首先定义清楚关系，然后才是执行步骤。

——Grace Murray Hopper，管理方法以及未来的计算机

第5章讨论了复制技术，即在不同节点上保存相同数据的多个副本。然而，面对一些海量数据集或非常高的查询压力，复制技术还不够，我们还需要将数据拆分成为分区，也称为分片<sup>注1</sup>。



## 术语澄清

这里我们所讨论的分区，在不同系统有着不同的称呼，例如它对应于MongoDB，Elasticsearch和SolrCloud中的*shard*，HBase的*region*，Bigtable中的*tablet*，Cassandra和Riak中的*vnode*，以及Couchbase中的*vBucket*。总体而言，分区是最普遍的术语。

分区通常是这样定义的，即每一条数据（或者每条记录，每行或每个文档）只属于某个特定分区。实现的方法有多种，稍后将逐一介绍。实际上，每个分区都可以视为一个完整的小型数据库，虽然数据库可能存在一些跨分区的操作。

采用数据分区的主要目的是提高可扩展性。不同的分区可以放在一个无共享集群（参阅第二部分关于无共享架构的定义）的不同节点上。这样一个大数据集可以分散在更

注1：本章所讨论的数据拆分指的是对数据集的划分，与网络分区问题（第8章讨论“一种节点间的网络故障”）完全不同。

多的磁盘上，查询负载也随之分布到更多的处理器上。

对单个分区进行查询时，每个节点对自己所在分区可以独立执行查询操作，因此添加更多的节点可以提高查询吞吐量。超大而复杂的查询尽管比较困难，但也可能做到跨节点的并行处理。

分区数据库最初在20世纪80年代由Teradata和Tandem NonStop SQL<sup>[1]</sup>等率先推出，最近又被一些NoSQL数据库和基于Hadoop的数据仓库重视起来。这些系统有些是为事务型负载设计的，有些是分析型（参阅第3章的“事务处理与分析处理”）。二者的差异会显著影响系统的优化策略，然而分区技术的基本原理则可以普遍适用。

本章我们将首先介绍切分大型数据集的若干方法，讨论数据索引如何影响分区。接下来讨论分区的再平衡，这对动态添加或删除节点非常重要。最后，我们将介绍数据库如何将请求路由到正确的分区并执行查询。

## 数据分区与数据复制

分区通常与复制结合使用，即每个分区在多个节点都存有副本。这意味着某条记录属于特定的分区，而同样的内容会保存在不同的节点上以提高系统的容错性。

一个节点上可能存储了多个分区。图6-1展示了主从复制模型与分区组合使用时数据的分布情况。由图可知，每个分区都有自己的主副本，例如被分配给某节点，而从副本则分配在其他一些节点。一个节点可能即是某些分区的主副本，同时又是其他分区的从副本。

第5章所讨论的所有复制相关的原理同样适用于对分区数据的复制。考虑到分区方案的选择通常独立于复制，因此本章将力求简洁，而省略与复制相关的内容。

## 键-值数据的分区

假设面临海量数据，现在需要切分它们，那么该如何决定哪些记录放在哪些节点上呢？

分区的主要目标是将数据和查询负载均匀分布在所有节点上。如果节点平均分担负载，那么理论上10个节点应该能够处理10倍的数据量和10倍于单个节点的读写吞吐量（忽略复制）。



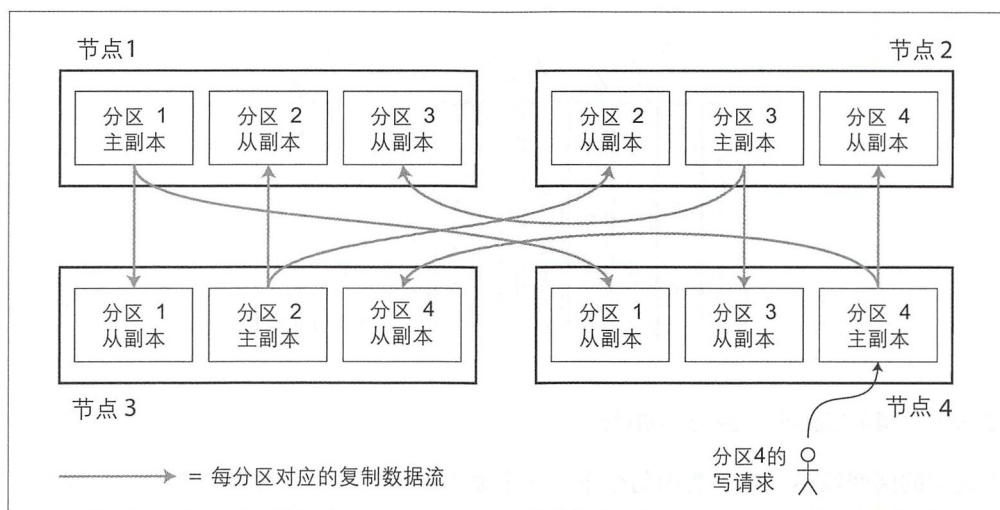


图6-1：组合使用复制和分区：每个节点同时充当某些分区的主副本和其他分区的从副本

而如果分区不均匀，则会出现某些分区节点比其他分区承担更多的数据量或查询负载，称之为倾斜。倾斜会导致分区效率严重下降，在极端情况下，所有的负载可能会集中在一个分区节点上，这就意味着10个节点9个空闲，系统的瓶颈在最繁忙的那个节点上。这种负载严重不成比例的分​​区即成为系统热点。

避免热点最简单的方法是将记录随机分配给所有节点上。这种方法可以比较均匀地分布数据，但是有一个很大的缺点：当试图读取特定的数据时，没有办法知道数据保存在哪个节点上，所以不得不并行查询所有节点。

可以改进上述方法。现在我们假设数据是简单的键-值数据模型，这意味着总是可以通过关键字来访问记录。例如，像一个纸质百科全书，可以通过标题来查找某一个条目；而所有的条目按字母序排序，因此可以做到快速查找条目。

## 基于关键字区间分区

一种分区方式是给每个分区分配一段连续的关键字或者关键字区间范围（以最小值和最大值来指示），如图6-2所示的纸质百科全书的卷目录。如果知道关键字区间的上下限，就可以轻松确定哪个分区包含这些关键字。如果还知道哪个分区分配在哪个节点，就可以直接向该节点发出请求（对于百科全书的例子，就是从书架上直接取到所要的书籍）。

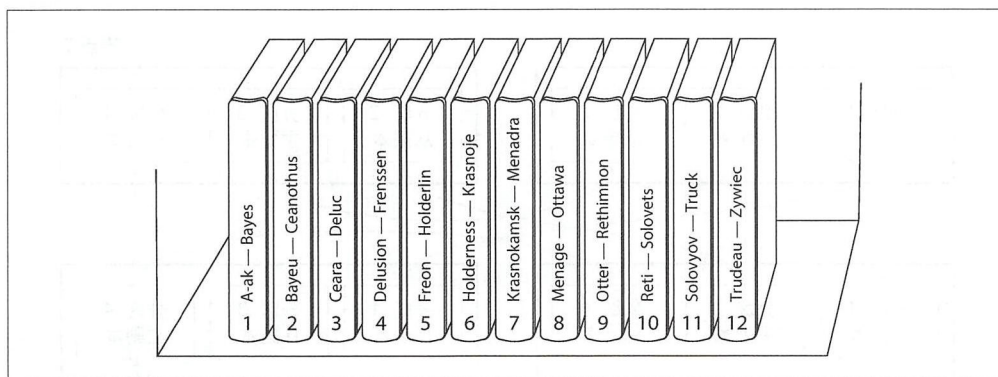


图6-2：百科全书按照关键字区间进行分区

关键字的区间段不一定非要均匀分布，这主要是因为数据本身可能就不均匀。例如，在图6-2中，卷1只包含以A和B开头的单词，但是卷12则包含了T、U、V、X、Y和Z开始的单词。如果只是简单地规定每个分区包含两个字母，则可能会导致一些卷比其他卷要大很多。为了更均匀地分布数据，分区边界理应适配数据本身的分布特征。

分区边界可以由管理员手动确定，或者由数据库自动选择（我们将在本章后面的“分区再平衡”中更详细地讨论）。采用这种分区策略的系统包括Bigtable，Bigtable的开源版本HBase<sup>[2,3]</sup>，RethinkDB和2.4版本之前MongoDB<sup>[4]</sup>。

每个分区内可以按照关键字排序保存（参阅第3章的“SSTables和LSM-Trees”）。这样可以轻松支持区间查询，即将关键字作为一个拼接起来的索引项从而一次查询得到多个相关记录（参阅第3章的“多列索引”）。例如，对于一个保存网络传感器数据的应用系统，选择测量的时间戳（年-月-日-时-分-秒）作为关键字，此时区间查询会非常有用，它可以快速获得某个月份内的所有数据。

然而，基于关键字的区间分区的缺点是某些访问模式会导致热点。如果关键字是时间戳，则分区对应于一个时间范围，例如每天一个分区。然而，当测量数据从传感器写入数据库时，所有的写入操作都集中在同一个分区（即当天的分区），这会导致该分区在写入时负载过高，而其他分区始终处于空闲状态<sup>[5]</sup>。

为了避免上述问题，需要使用时间戳以外的其他内容作为关键字的第一项。例如，可以在时间戳前面加上传感器名称作为前缀，这样首先由传感器名称，然后按时间进行分区。假设同时有许多传感器处于活动状态，则写入负载最终会比较均匀地分布在多个节点上。接下来，当需要获取一个时间范围内、多个传感器的数据时，可以根据传感器名称，各自执行区间查询。

## 基于关键字哈希值分区

对于上述数据倾斜与热点问题，许多分布式系统采用了基于关键字哈希函数的方式来分区。

一个好的哈希函数可以处理数据倾斜并使其均匀分布。例如一个处理字符串的32位哈希函数，当输入某个字符串，它会返回一个0和 $2^{32}-1$ 之间近似随机分布的数值。即使输入的字符串非常相似，返回的哈希值也会在上述数字范围内均匀分布。

用于数据分区目的的哈希函数不需要在加密方面很强：例如，Cassandra和MongoDB使用MD5，Voldemort使用Fowler-Noll-Vo函数。许多编程语言也有内置的简单哈希函数（主要用于哈希表），但是要注意这些内置的哈希函数可能并不适合分区，例如，Java的Object.hashCode和Ruby的Object#hash，同一个键在不同的进程中可能返回不同的哈希值<sup>[6]</sup>。

一旦找到合适的关键字哈希函数，就可以为每个分区分配一个哈希范围（而不是直接作用于关键字范围），关键字根据其哈希值的范围划分到不同的分区中。如图6-3所示。

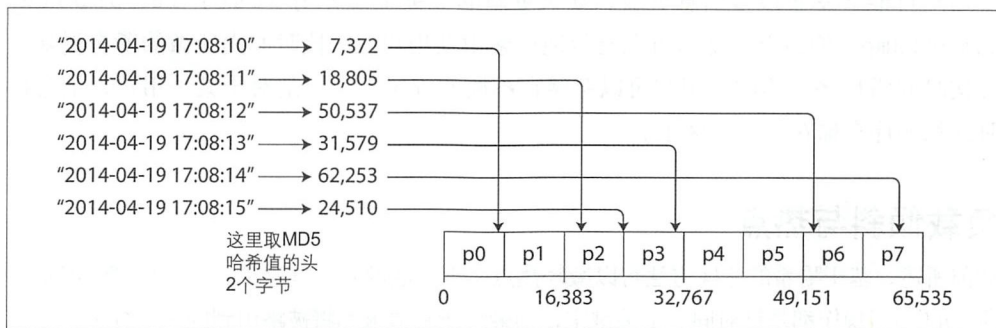


图6-3：基于关键字的哈希值进行分区

这种方法可以很好地将关键字均匀地分配到多个分区中。分区边界可以是均匀间隔，也可以是伪随机选择（在这种情况下，该技术有时被称为一致性哈希）。

然而，通过关键字哈希进行分区，我们丧失了良好的区间查询特性。即使关键字相邻，但经过哈希之后会分散在不同的分区中，区间查询就失去了原有的有序相邻的特性。在MongoDB中，如果启用了基于哈希的分片模式，则区间查询会发送到所有的分区上<sup>[4]</sup>，而Riak<sup>[9]</sup>、Couchbase<sup>[10]</sup>和Voldemort干脆就不支持关键字上的区间查询。



## 一致性哈希

一致性哈希，由Karger等人首先提出<sup>[7]</sup>，是一种平均分配负载的方法，最初用于内容分发网络（CDN）等互联网缓存系统。它采用随机选择的分区边界来规避中央控制或分布式共识。请注意，此处的一致性与副本一致性（第5章）或ACID一致性（第7章）没有任何关联，它只描述了数据动态平衡的一种方法。

正如后面“分区再平衡”一节将要介绍的，这种特殊的分区方法对于数据库实际效果并不是很好，所以目前很少使用（虽然某些数据库的文档仍采用一致性哈希的术语，但其实并不准确）。为避免混淆，我们此处只用术语哈希分区。

Cassandra则在两种分区策略之间做了一个折中<sup>[11-13]</sup>。Cassandra中的表可以声明为由多个列组成的复合主键。复合主键只有第一部分可用于哈希分区，而其他列则用作组合索引来对Cassandra SSTable中的数据进行排序。因此，它不支持在第一列上进行区间查询，但如果为第一列指定好了固定值，可以对其他列执行高效的区间查询。

组合索引为一对多的关系提供了一个优雅的数据模型。例如，在社交网站上，一个用户可能会发布很多消息更新。如果更新的关键字设置为（`user_id`, `update_timestamp`）的组合，那么可以有效地检索由某用户在一段时间内所做的所有更新，且按时间戳排序。不同的用户可以存储在不同的分区上，但是对于某一用户，消息按时间戳顺序存储在一个分区上。

## 负载倾斜与热点

如前所述，基于哈希的分区方法可以减轻热点，但无法做到完全避免。一个极端情况是，所有的读/写操作都是针对同一个关键字，则最终所有请求都将被路由到同一个分区。

这种负载或许并不普遍，但也并非不可能：例如，社交媒体网站上，一些名人用户有数百万的粉丝，当其发布一些热点事件时可能会引发一场访问风暴<sup>[14]</sup>，出现大量的对相同关键字的写操作（其中关键字可能是名人的用户ID，或者人们正在评论的事件ID）。此时，哈希起不到任何帮助作用，因为两个相同ID的哈希值仍然相同。

大多数的系统今天仍然无法自动消除这种高度倾斜的负载，而只能通过应用层来减轻倾斜程度。例如，如果某个关键字被确认为热点，一个简单的技术就是在关键字的开头或结尾处添加一个随机数。只需一个两位数的十进制随机数就可以将关键字的写操作分布到100个不同的关键字上，从而分配到不同的分区上。

但是，随之而来的问题是，之后的任何读取都需要些额外的工作，必须从所有100个

关键字中读取数据然后进行合并。因此通常只对少量的热点关键字附加随机数才有意义；而对于写入吞吐量低的绝大多数关键字，这些都意味着不必要的开销。此外，还需要额外的元数据来标记哪些关键字进行了特殊处理。

也许将来某一天，数据系统能够自动检测负载倾斜情况，然后自动处理这些倾斜的负载。但截至目前，仍然需要开发者自己结合应用来综合权衡。

## 分区与二级索引

我们之前所讨论的分区方案都依赖于键-值数据模型。键-值模型相对简单，即都是通过关键字来访问记录，自然可以根据关键字来确定分区，并将读写请求路由到负责该关键字的分区上。

但是，如果涉及二级索引，情况会变得复杂（参阅第3章的“其他索引结构”）。二级索引通常不能唯一标识一条记录，而是用来加速特定值的查询，例如查找用户123的所有操作，找到所有含有hogwash的文章，查找所有颜色为红色的汽车等。

二级索引是关系数据库的必备特性，在文档数据库中应用也非常普遍。但考虑到其复杂性，许多键-值存储（如HBase和Voldemort）并不支持二级索引；但其他一些如Riak则开始增加对二级索引的支持。此外，二级索引技术也是Solr和Elasticsearch等全文索引服务器存在之根本。

二级索引带来的主要挑战是它们不能规整地映射到分区中。有两种主要的方法来支持对二级索引进行分区：基于文档的分区和基于词条的分区。

### 基于文档分区的二级索引

假设一个销售二手车的网站（见图6-4）。每个列表都有一个唯一的文档ID，用此ID对数据库进行分区，例如，ID 0到499归分区0，ID 500到999划为分区1。

现在用户需要搜索汽车，可以持按汽车颜色和厂商进行过滤，所以需要在颜色和制造商上设定二级索引（在文档数据库中这些都是字段；在关系数据库中则是列）。声明这些索引之后，数据库会自动<sup>注2</sup>创建索引。例如，每当一辆红色汽车添加到数据库中，数据库分区会自动将其添加到索引条目为“color:red”的文档ID列表中。

---

注2：如果数据库仅支持键-值模型，则可能需要在应用层代码中创建从值到文档ID的映射来实现二级索引。如果采用这种方式，需要格外小心，确保索引与原数据库系统保持数据一致。发生任何竞争条件以及中间写失败（出现只保存了部分修改的情况）都很容易导致数据不同步。请参阅第7章“多对象事务”。

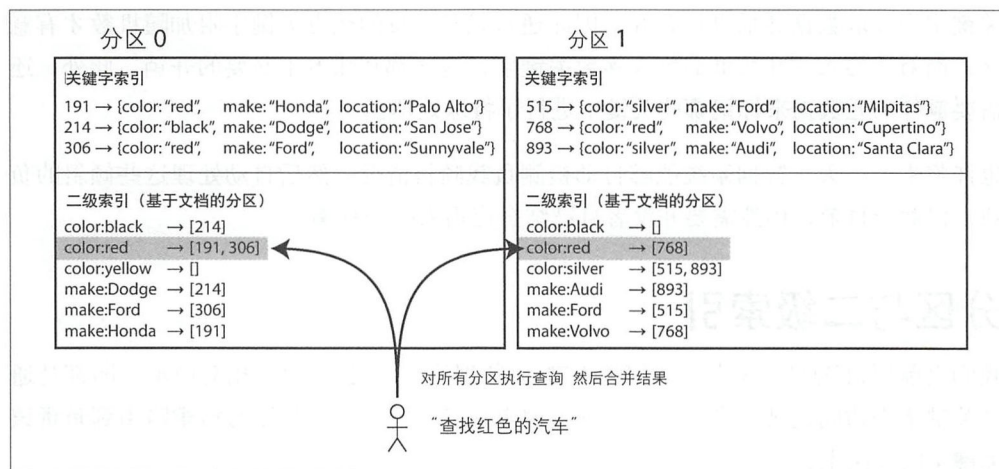


图6-4：基于文档的二级索引分区

在这种索引方法中，每个分区完全独立，各自维护自己的二级索引，且只负责自己分区内的文档而不关心其他分区中数据。每当需要写数据库时，包括添加，删除或更新文档等，只需要处理包含目标文档ID的那一个分区。因此文档分区索引也被称为本地索引，而不是全局索引，后者将在本章后面介绍。

但读取时需要注意：除非对文档ID做了特别的处理，否则不太可能所有特定颜色或特定品牌的汽车都放在一个分区中，例如图6-4中，红色汽车就出现在分区0和分区1中。因此，如果想要搜索红色汽车，就需要将查询发送到所有的分区，然后合并所有返回的结果。

这种查询分区数据库的方法有时也称为分散/聚集，显然这种二级索引的查询代价高昂。即使采用了并行查询，也容易导致读延迟显著放大（参阅第1章的“实践中的百分位数”）。尽管如此，它还是广泛用于实践：MongoDB、Riak<sup>[15]</sup>、Cassandra<sup>[16]</sup>、Elasticsearch<sup>[17]</sup>、SolrCloud<sup>[18]</sup>和VoltDB<sup>[19]</sup>都支持基于文档分区二级索引。大多数数据库供应商都建议用户自己来构建合适的分区方案，尽量由单个分区满足二级索引查询，但现实往往难以如愿，尤其是当查询中可能引用多个二级索引时（例如同时指定颜色和制造商两个条件）。

## 基于词条的二级索引分区

另一种方法，我们可以对所有的数据构建全局索引，而不是每个分区维护自己的本地索引。而且，为避免成为瓶颈，不能将全局索引存储在一个节点上，否则就破坏了设



计分区均衡的目标。所以，全局索引也必须进行分区，且可以与数据关键字采用不同的分区策略。

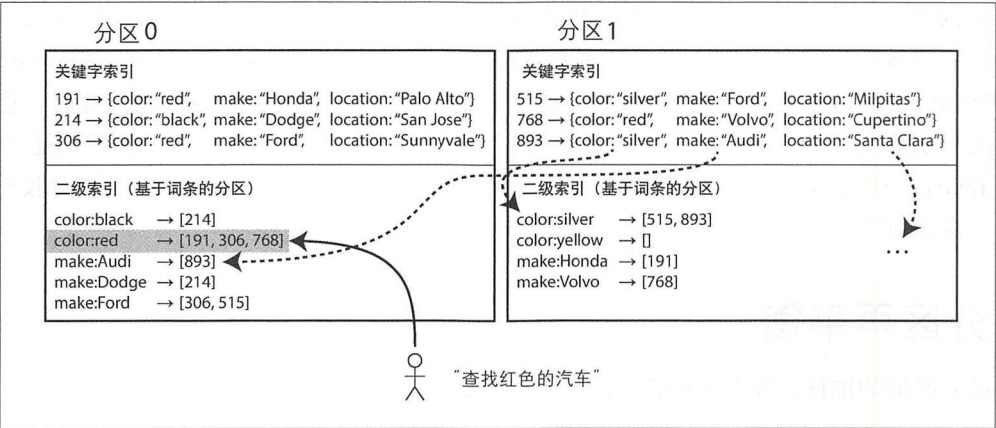


图6-5：基于词条的二级索引分区

以图6-5为例：所有数据分区中的颜色为红色的汽车被收录到在索引color:red中，而索引本身也是分区的，例如从a到r开始的颜色放在分区0中，从s到z的颜色放在分区1中。类似的，汽车制造商的索引也被分区（两个分区的边界分别是字母f和字母h）。

我们将这种索引方案称为词条分区，它以待查找的关键字本身作为索引。例如颜色color:red。名字词条源于全文索引（一种特定类型的二级索引），term指的是文档中出现的所有单词的集合。

和前面讨论的方法一样，可以直接通过关键词来全局划分索引，或者对其取哈希值。直接分区的好处是可以支持高效的区间查询（例如，查询汽车报价在某个值以上）；而采用哈希的方式则可以更均匀的划分分区。

这种全局的词条分区相比于文档分区索引的主要优点是，它的读取更为高效，即它不需要采用scatter/gather对所有的分区都执行一遍查询，相反，客户端只需要向包含词条的那一个分区发出读请求。然而全局索引的不利之处在于，写入速度较慢且非常复杂，主要因为单个文档的更新时，里面可能会涉及多个二级索引，而二级索引的分区又可能完全不同甚至在不同的节点上，由此势必引入显著的写放大。

理想情况下，索引应该时刻保持最新，即写入的数据要立即反映在最新的索引上。但是，对于词条分区来讲，这需要一个跨多个相关分区的分布式事务支持，写入速度会

受到极大的影响，所以现有的数据库都不支持同步更新二级索引（参阅第7章和第9章）。

实践中，对全局二级索引的更新往往都是异步的（也就意味着，如果在写入之后马上去读索引，那么刚刚发生的更新可能还没有反映在索引中）。例如，Amazon DynamoDB的二级索引通常可以在1秒之内完成更新，但当底层设施出现故障时，也有可能需要等待很长的时间<sup>[20]</sup>。其他使用全局索引的系统还包括Riak的搜索功能<sup>[21]</sup>和Oracle数据仓库，后者允许用户来选择是使用本地还是全局索引<sup>[22]</sup>。在第12章我们会重新讨论如何实现全局二级索引。

## 分区再平衡

随着时间的推移，数据库可能总会出现某些变化：

- 查询压力增加，因此需要更多的CPU来处理负载。
- 数据规模增加，因此需要更多的磁盘和内存来存储数据。
- 节点可能出现故障，因此需要其他机器来接管失效的节点。

所有这些变化都要求数据和请求可以从一个节点转移到另一个节点。这样一个迁移负载的过程称为再平衡（或者动态平衡）。无论对于哪种分区方案，分区再平衡通常至少要满足：

- 平衡之后，负载、数据存储、读写请求等应该在集群范围更均匀地分布。
- 再平衡执行过程中，数据库应该可以继续正常提供读写服务。
- 避免不必要的负载迁移，以加快动态再平衡，并尽量减少网络和磁盘I/O影响。

## 动态再平衡的策略

将分区对应到节点上存在多种不同的分配策略，这里逐一介绍：

### 为什么不用取模？

我们在前面提到（见图6-3），最好将哈希值划分为不同的区间范围，然后将每个区间分配给一个分区。例如，区间 $[0, b_0)$ 对应于分区0， $[b_0, b_1)$ 对应分区1等。

也许你会问为什么不直接使用 $mod$ （许多编程语言里的取模运算符 $\%$ ）。例如， $hash(key) \bmod 10$ 会返回一个介于0和9之间的数字，如果有10个节点，则依次对应节点0到9，这似乎是将每个关键字分配到节点的最简单方法。

对节点数取模方法的问题是，如果节点数N发生了变化，会导致很多关键字需要从现有的节点迁移到另一个节点。例如，假设 $hash(key) = 123456$ ，假定最初是10个节点，那么这个关键字应该放在节点6 ( $123456 \bmod 10 = 6$ )；当节点数增加到11时，它需要移动到节点3 ( $123456 \bmod 11 = 3$ )；当继续增长到12个节点时，又需要移动到节点0 ( $123456 \bmod 12 = 0$ )。这种频繁的迁移操作大大增加了再平衡的成本。

因此我们需要一种减少迁移数据的方法。

## 固定数量的分区

幸运的是，有一个相当简单的解决方案：首先，创建远超实际节点数的分区数，然后为每个节点分配多个分区。例如，对于一个10节点的集群，数据库可以从一开始就逻辑划分为1000个分区，这样大约每个节点承担100个分区。

接下来，如果集群中添加了一个新节点，该新节点可以从每个现有的节点上匀走几个分区，直到分区再次达到全局平衡。该过程如图6-6所示。如果从集群中删除节点，则采取相反的均衡措施。

选中的整个分区会在节点之间迁移，但分区的总数量仍维持不变，也不会改变关键字到分区的映射关系。这里唯一要调整的是分区与节点的对应关系。考虑到节点间通过网络传输数据总是需要些时间，这样调整可以逐步完成，在此期间，旧的分区仍然可以接收读写请求。

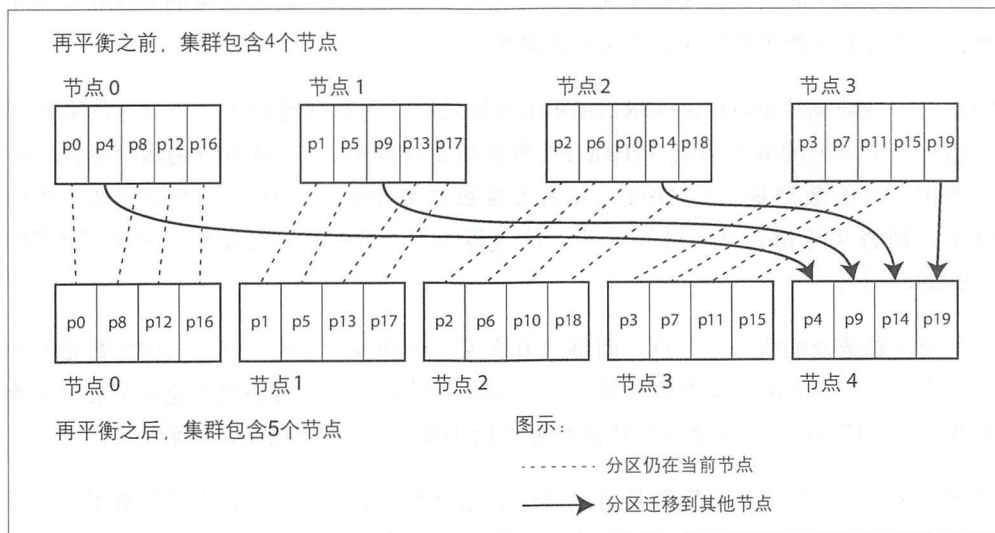


图6-6：新节点加入到集群，每个节点上承担更多的分区数据



原则上，也可以将集群中的不同的硬件配置因素考虑进来，即性能更强大的节点将分配更多的分区，从而分担更多的负载。

目前，Riak<sup>[15]</sup>、Elasticsearch<sup>[24]</sup>、Couchbase<sup>[10]</sup>和Voldemort<sup>[25]</sup>都支持这种动态平衡方法。

使用该策略时，分区的数量往往在数据库创建时就确定好，之后不会改变。原则上也可以拆分和合并分区（稍后介绍），但固定数量的分区使得相关操作非常简单，因此许多采用固定分区策略的数据库决定不支持分区拆分功能。所以，在初始化时，已经充分考虑将来扩容增长的需求（未来可能拥有的最大节点数），设置一个足够大的分区数。而每个分区也有些额外的管理开销，选择过高的数字可能会有副作用。

如果数据集的总规模高度不确定或可变（例如，开始非常小，但随着时间的推移可能会变得异常庞大），此时如何选择合适的分区数就有些困难。每个分区包含的数据量的上限是固定的，实际大小应该与集群中的数据总量成正比。如果分区里的数据量非常大，则每次再平衡和节点故障恢复的代价就很大；但是如果一个分区太小，就会产生太多的开销。分区大小应该“恰到好处”，不要太大，也不能过小，如果分区数量固定了但总数据量却高度不确定，就难以达到一个最佳取舍点。

## 动态分区

对于采用关键字区间分区的数据库，如果边界设置有问题，最终可能会出现所有数据都挤在一个分区而其他分区基本为空，那么设定固定边界、固定数量的分区将非常不便：而手动去重新配置分区边界又非常繁琐。

因此，一些数据库如HBase和RethinkDB等采用了动态创建分区。当分区的数据增长超过一个可配的参数阈值（HBase上默认值是10GB），它就拆分为两个分区，每个承担一半的数据量<sup>[26]</sup>。相反，如果大量数据被删除，并且分区缩小到某个阈值以下，则将其与相邻分区进行合并。该过程类似于B树的分裂操作（参阅第3章的“B-tree”）。

每个分区总是分配给一个节点，而每个节点可以承载多个分区，这点与固定数量的分区一样。当一个大的分区发生分裂之后，可以将其中的一半转移到其他某节点以平衡负载。对于HBase，分区文件的传输需要借助HDFS（底层分布式文件系统）<sup>[3]</sup>。

动态分区的一个优点是分区数量可以自动适配数据总量。如果只有少量的数据，少量的分区就足够了，这样系统开销很小；如果有大量的数据，每个分区的大小则被限制在一个可配的最大值<sup>[23]</sup>。

但是，需要注意的是，对于一个空的数据库，因为没有任何先验知识可以帮助确定分区的边界，所以会从一个分区开始。可能数据集很小，但直到达到第一个分裂点之前，所有的写入操作都必须由单个节点来处理，而其他节点则处于空闲状态。为了缓解这个问题，HBase和MongoDB允许在一个空的数据库上配置一组初始分区（这被称为预分裂）。对于关键字区间分区，预分裂要求已经知道一些关键字的分布情况<sup>[4,26]</sup>。

动态分区不仅适用于关键字区间分区，也适用于基于哈希的分区策略。MongoDB从版本2.4开始，同时支持二者，并且都可以动态分裂分区。

## 按节点比例分区

采用动态分区策略，拆分和合并操作使每个分区的大小维持在设定的最小值和最大值之间，因此分区的数量与数据集的大小成正比关系。另一方面，对于固定数量的分区方式，其每个分区的大小也与数据集的大小成正比。两种情况，分区的数量都与节点数无关。

Cassandra和Ketama则采用了第三种方式，使分区数与集群节点数成正比关系。换句话说，每个节点具有固定数量的分区<sup>[23,27,28]</sup>。此时，当节点数不变时，每个分区的大小与数据集大小保持正比的增长关系；当节点数增加时，分区则会调整变得更小。较大的数据量通常需要大量的节点来存储，因此这种方法也使每个分区大小保持稳定。

当一个新节点加入集群时，它随机选择固定数量的现有分区进行分裂，然后拿走这些分区的一半数据量，将另一半数据留在原节点。随机选择可能会带来不太公平的分区分裂，但是当平均分区数量较大时（Cassandra默认情况下，每个节点有256个分区），新节点最终会从现有节点中拿走相当数量的负载。Cassandra 在3.0时推出了改进算法，可以避免上述不公平的分裂<sup>[29]</sup>。

随机选择分区边界的前提要求采用基于哈希分区（可以从哈希函数产生的数字范围里设置边界）。这种方法也最符合本章开头所定义一致性哈希<sup>[7]</sup>。一些新设计的哈希函数也可以以较低的元数据开销达到类似的效果<sup>[8]</sup>。

## 自动与手动再平衡操作

动态平衡另一个重要问题我们还没有考虑：它是自动执行还是手动方式执行？

全自动式的再平衡（即由系统自动决定何时将分区从一个节点迁移到另一个节点，不需要任何管理员的介入）与纯手动方式（即分区到节点的映射由管理员来显式配置）

之间，可能还有一个过渡阶段。例如，Couchbase，Riak和Voldemort会自动生成一个分区分区的建议方案，但需要管理员的确认才能生效。

全自动式再平衡会更加方便，它在正常维护之外所增加的操作很少。但是，也有可能出现结果难以预测的情况。再平衡总体讲是个比较昂贵的操作，它需要重新路由请求并将大量数据从一个节点迁移到另一个节点。万一执行过程中出现异常，会使网络或节点的负载过重，并影响其他请求的性能。

将自动平衡与自动故障检测相结合也可能存在一些风险。例如，假设某个节点负载过重，对请求的响应暂时受到影响，而其他节点可能会得到结论：该节点已经失效；接下来激活自动平衡来转移其负载。客观上这会加重该节点、其他节点以及网络的负荷，可能会使总体情况变得更糟，甚至导致级联式的失效扩散。

出于这样的考虑，让管理员介入到再平衡可能是个更好的选择。它的确比全自动过程响应慢一些，但它可以有效防止意外发生。

## 请求路由

现在我们已经将数据集分布到多个节点上，但是仍然有一个悬而未决的问题：当客户端需要发送请求时，如何知道应该连接哪个节点？如果发生了分区再平衡，分区与节点的对应关系随之还会变化。为了回答该问题，我们需要一段处理逻辑来感知这些变化，并负责处理客户端的连接，例如想要读/写关键字“foo”，需要连接哪个IP地址和哪个端口号。

这其实属于一类典型的服务发现问题，服务发现并不限于数据库，任何通过网络访问的系统都有这样的需求，尤其是当服务目标支持高可用时（在多台机器上有冗余配置）。许多公司已经开发了自己的内部服务发现工具，其中很多已经开源<sup>[30]</sup>。

概括来讲，这个问题有以下几种不同的处理策略（分别如图6-7所示的三种情况）：

1. 允许客户端链接任意的节点（例如，采用循环式的负载均衡器）。如果某节点恰好拥有所请求的分区，则直接处理该请求；否则，将请求转发到下一个合适的节点，接收答复，并将答复返回给客户端。
2. 将所有客户端的请求都发送到一个路由层，由后者负责将请求转发到对应的分区节点上。路由层本身不处理任何请求，它仅充一个分区感知的负载均衡器。
3. 客户端感知分区和节点分配关系。此时，客户端可以直接连接到目标节点，而不需要任何中介。



不管哪种方法，核心问题是：作出路由决策的组件（可能是某个节点，路由层或客户端）如何知道分区与节点的对应关系以及其变化情况？

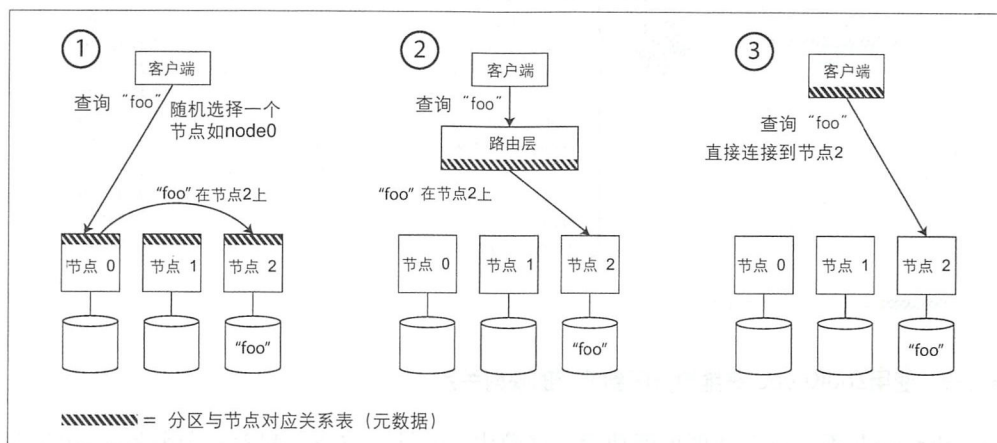


图6-7：将请求路由到正确分区节点的三种方式

这其实是一个很有挑战性的问题，所有参与者都要达成共识这一点很重要。否则请求可能被发送到错误的节点，而没有得到正确处理。分布式系统中有专门的共识协议算法，但通常难以正确实现（详见第9章）。

许多分布式数据系统依靠独立的协调服务（如ZooKeeper）跟踪集群范围内的元数据，如图6-8所示。每个节点都向ZooKeeper中注册自己，ZooKeeper维护了分区到节点的最终映射关系。其他参与者（如路由层或分区感知的客户端）可以向ZooKeeper订阅此信息。一旦分区发生了改变，或者添加、删除节点，ZooKeeper就会主动通知路由层，这样使路由信息保持最新状态。

例如，LinkedIn的Espresso使用Helix<sup>[31]</sup>进行集群管理（底层是ZooKeeper），实现了图6-8所示的请求路由层。HBase，SolrCloud和Kafka也使用ZooKeeper来跟踪分区分配情况。MongoDB有类似的设计，但它依赖于自己的配置服务器和mongos守护进程来充当路由层。

Cassandra和Riak则采用了不同的方法，它们在节点之间使用gossip协议来同步群集状态的变化。请求可以发送到任何节点，由该节点负责将其转发到目标分区节点（图6-7中的方法1）。这种方式增加了数据库节点的复杂性，但是避免了对ZooKeeper之类的外部协调服务的依赖。

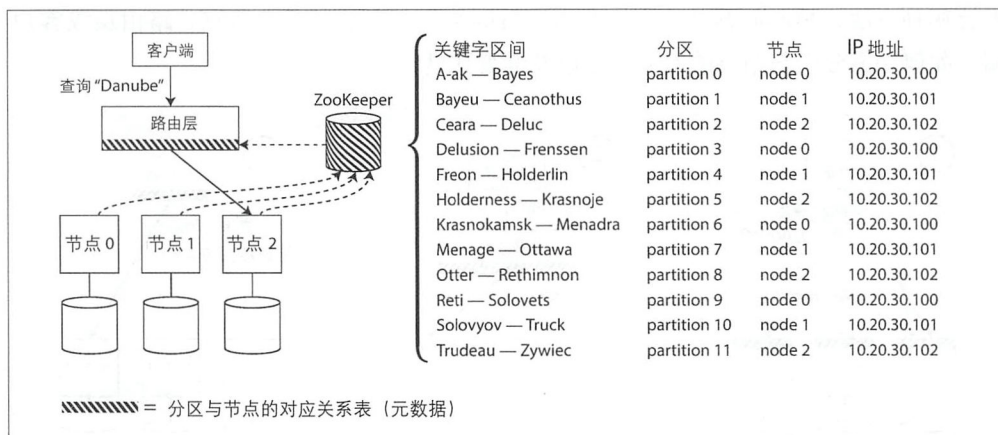


图 6-8：使用ZooKeeper来维护分区到节点的映射关系

Couchbase并不支持自动再平衡功能，这简化了设计。它通过配置一个名为*moxi*的路由选择层，向集群节点学习最新的路由变化<sup>[32]</sup>。

当使用路由层或随机选择节点发送请求时，客户端仍然需要知道目标节点的IP地址。IP地址的变化往往没有分区-节点变化那么频繁，采用DNS通常就足够了。

## 并行查询执行

到目前为止，我们只关注了读取或写入单个关键字这样简单的查询（对于文档分区的二级索引，里面要求分散/聚集查询）。这基本上也是大多数NoSQL分布式数据存储所支持的访问类型。

然而对于大规模并行处理（massively parallel processing, MPP）这一类主要用于数据分析的关系数据库，在查询类型方面要复杂得多。典型的数据仓库查询包含多个联合、过滤、分组和聚合操作。MPP查询优化器会将复杂的查询分解成许多执行阶段和分区，以便在集群的不同节点上并行执行。尤其是涉及全表扫描这样的查询操作，可以通过并行执行获益颇多。

数据仓库中快速并行执行查询可以作为单独的话题。考虑到分析业务的重要性，目前它已得到了广泛的商业关注。我们将在第10章中讨论并行查询执行所需的一些技术。有关并行数据库更多相关技术细节请参考文献<sup>[1,33]</sup>。

## 小结

本章，我们探讨了将大规模数据集划分成更小子集的多种方法。数据量如果太大，单

台机器进行存储和处理就会成为瓶颈，因此需要引入数据分区机制。分区的目地是通过多台机器均匀分布数据和查询负载，避免出现热点。这需要选择合适的数据库分区方案，在节点添加或删除时重新动态平衡分区。

我们讨论了两种主要的分区方法：

- 基于关键字区间的分区。先对关键字进行排序，每个分区只负责一段包含最小到最大关键字范围的一段关键字。对关键字排序的优点是可以支持高效的区间查询，但是如果应用程序经常访问与排序一致的某段关键字，就会存在热点的风险。采用这种方法，当分区太大时，通常将其分裂为两个子区间，从而动态地再平衡分区。
- 哈希分区。将哈希函数作用于每个关键字，每个分区负责一定范围的哈希值。这种方法打破了原关键字的顺序关系，它的区间查询效率比较低，但可以更均匀地分配负载。采用哈希分区时，通常事先创建好足够多（但固定数量）的分区，让每个节点承担多个分区，当添加或删除节点时将某些分区从一个节点迁移到另一个节点，也可以支持动态分区。

混合上述两种基本方法也是可行的，例如使用复合键：键的一部分来标识分区，而另一部分来记录排序后的顺序。

我们还讨论了分区与二级索引，二级索引也需要进行分区，有两种方法：

- 基于文档来分区二级索引（本地索引）。二级索引存储在与关键字相同的分区中，这意味着写入时我们只需要更新一个分区，但缺点是读取二级索引时需要在所有分区上执行scatter/gather。
- 基于词条来分区二级索引（全局索引）。它是基于索引的值而进行的独立分区。二级索引中的条目可能包含来自关键字的多个分区里的记录。在写入时，不得不更新二级索引的多个分区；但读取时，则可以从单个分区直接快速提取数据。

最后，我们讨论了如何将查询请求路由到正确的分区，包括简单的分区感知负载均衡器，以及复杂的并行查询执行引擎。

理论上，每个分区基本保持独立运行，这也是为什么我们试图将分区数据库分布、扩展到多台机器上。但是，如果写入需要跨多个分区，情况就会格外复杂，例如，如果其中一个分区写入成功，但另一个发生失败，接下来会发生什么？我们将在下面的章节中讨论类似这样的技术挑战。



## 参考文献

- [1] David J. DeWitt and Jim N. Gray: “Parallel Database Systems: The Future of High Performance Database Systems,” *Communications of the ACM*, volume 35, number 6, pages 85-98, June 1992. doi:10.1145/129888.129894.
- [2] Lars George: “HBase vs. BigTable Comparison,” *larsgeorge.com*, November 2009.
- [3] “The Apache HBase Reference Guide,” Apache Software Foundation, *hbase.apache.org*, 2014.
- [4] MongoDB, Inc.: “New Hash-Based Sharding Feature in MongoDB 2.4,” *blog.mongodb.org*, April 10, 2013.
- [5] Ikai Lan: “App Engine Datastore Tip: Monotonically Increasing Values Are Bad,” *ikaisays.com*, January 25, 2011.
- [6] Martin Kleppmann: “Java’s hashCode Is Not Safe for Distributed Systems,” *martin.kleppmann.com*, June 18, 2012.
- [7] David Karger, Eric Lehman, Tom Leighton, et al.: “Consistent Hashing and Random Trees: Distributed Caching Protocols for Relieving Hot Spots on the World Wide Web,” at *29th Annual ACM Symposium on Theory of Computing (STOC)*, pages 654-663, 1997. doi:10.1145/258533.258660.
- [8] John Lamping and Eric Veach: “A Fast, Minimal Memory, Consistent Hash Algorithm,” *arxiv.org*, June 2014.
- [9] Eric Redmond: “A Little Riak Book,” Version 1.4.0, Basho Technologies, September 2013.
- [10] “Couchbase 2.5 Administrator Guide,” Couchbase, Inc., 2014.
- [11] Avinash Lakshman and Prashant Malik: “Cassandra — A Decentralized Structured Storage System,” at *3rd ACM SIGOPS International Workshop on Large Scale Distributed Systems and Middleware (LADIS)*, October 2009.
- [12] Jonathan Ellis: “Facebook’s Cassandra Paper, Annotated and Compared to Apache Cassandra 2.0,” *datastax.com*, September 12, 2013.
- [13] “Introduction to Cassandra Query Language,” DataStax, Inc., 2014.

- [14] Samuel Axon: “3% of Twitter’s Servers Dedicated to Justin Bieber,” *mashable.com*, September 7, 2010.
- [15] “Riak 1.4.8 Docs,” Basho Technologies, Inc., 2014.
- [16] Richard Low: “The Sweet Spot for Cassandra Secondary Indexing,” *wentnet.com*, October 21, 2013.
- [17] Zachary Tong: “Customizing Your Document Routing,” *elasticsearch.org*, June 3, 2013.
- [18] “Apache Solr Reference Guide,” Apache Software Foundation, 2014.
- [19] Andrew Pavlo: “H-Store Frequently Asked Questions,” *hstore.cs.brown.edu*, October 2013.
- [20] “Amazon DynamoDB Developer Guide,” Amazon Web Services, Inc., 2014.
- [21] Rusty Klophaus: “Difference Between 2I and Search,” email to *riak-users* mailinglist, *lists.basho.com*, October 25, 2011.
- [22] Donald K. Burleson: “Object Partitioning in Oracle,” *dba-oracle.com*, November 8, 2000.
- [23] Eric Evans: “Rethinking Topology in Cassandra,” at *ApacheCon Europe*, November 2012.
- [24] Rafał Kuć: “Reroute API Explained,” *elasticsearchserverbook.com*, September 30, 2013.
- [25] “Project Voldemort Documentation,” *project-voldemort.com*.
- [26] Enis Soztutar: “Apache HBase Region Splitting and Merging,” *hortonworks.com*, February 1, 2013.
- [27] Brandon Williams: “Virtual Nodes in Cassandra 1.2,” *datastax.com*, December 4, 2012.
- [28] Richard Jones: “libketama: Consistent Hashing Library for Memcached Clients,” *metabrew.com*, April 10, 2007.

- [29] Branimir Lambov: “New Token Allocation Algorithm in Cassandra 3.0,” *datastax.com*, January 28, 2016.
- [30] Jason Wilder: “Open-Source Service Discovery,” *jasonwilder.com*, February 2014.
- [31] Kishore Gopalakrishna, Shi Lu, Zhen Zhang, et al.: “Untangling Cluster Management with Helix,” at *ACM Symposium on Cloud Computing (SoCC)*, October 2012. doi:10.1145/2391229.2391248.
- [32] “Moxi 1.8 Manual,” Couchbase, Inc., 2014.
- [33] Shivnath Babu and Herodotos Herodotou: “Massively Parallel Databases and MapReduce Systems,” *Foundations and Trends in Databases*, volume 5, number 1, pages 1-104, November 2013. doi:10.1561/19000000036.





有些人抱怨，常用的两阶段提交在性能和可用性方面代价太高。而我们认为事务滥用和过度使用所引入的性能瓶颈应该主要由应用层来解决，而不是简单的抛弃事务。

——James Corbett 等，Spanner: 来自Google的全球分布式数据库（OSDI, 2012）

在一个苛刻的数据存储环境中，会有许多可能出错的情况，例如：

- 数据库软件或硬件可能会随时失效（包括正在执行写操作的过程中）。
- 应用程序可能随时崩溃（包括一系列操作执行到中间某一步）。
- 应用与数据库节点之间的链接可能随时会中断，数据库节点之间也存在同样问题。
- 多个客户端可能同时写入数据库，导致数据覆盖。
- 客户端可能读到一些无意义的、部分更新的数据。
- 客户端之间由于边界条件竞争所引入的各种奇怪问题。

为了系统高可靠的目标，我们必须处理好上述问题，万一发生类似情况确保不会导致系统级的失效。然而，完善的容错机制需要大量的工作，要仔细考虑各种可能出错的可能，并进行充分的测试才能确保方案切实可靠。

近十年来，事务技术一直是简化这些问题的首选机制。事务将应用程序的多个读、写操作捆绑在一起成为一个逻辑操作单元。即事务中的所有读写是一个执行的整体，整

个事务要么成功（提交）、要么失败（中止或回滚）。如果失败，应用程序可以安全地重试。这样，由于不需要担心部分失败的情况（无论出于何种原因），应用层的错误处理就变得简单很多。

也许对于那些浸淫此领域多年的读者来说，事务的概念如此简单，但细究起来或许并非如此。事务不是一个天然存在的东西，它是被人为创造出来，目的是简化应用层的编程模型。有了事务，应用程序可以不用考虑某些内部潜在的错误以及复杂的并发性问题，这些都可以交给数据库来负责处理（我们称之为安全性保证）。

然而并非每个应用程序都需要事务机制，有时可以弱化事务处理或完全放弃事务（例如，为了实现更高的性能或更高的可用性）。一些安全相关的属性也可能会避免引入事务。

那该如何判断是否需要事务呢？为了回答这个问题，我们首先需要确切地理解事务能够提供哪些安全性保证，背后的代价又是什么。事务的概念看似简单，实际上却有許多微妙而关键的细节值得研究。

本章，我们将首先分析可能出错的各种场景，探讨数据库防范这些问题的基本方法和算法设计。特别是在并发控制方面，深入讨论可能各种竞争条件，数据库所提供的多种隔离级别，例如读提交、快照隔离和可串行化等。

本章的内容可适用于单节点和分布式场景。在第8章中，我们将重点讨论分布式系统的特殊挑战。

## 深入理解事务

目前几乎所有的关系数据库和一些非关系数据库都支持事务处理。它们大多数都沿用了和IBM于1975年推出的第一个SQL数据库System R<sup>[1-3]</sup>相似的总体设计。尽管在一些具体实现方面有些不同，但事务的概念在这四十年中几乎没有发生变化，换句话说，MySQL、PostgreSQL、Oracle、SQL Server等系统实现的事务与当年System R非常相似。

然而21世纪末，非关系（NoSQL）数据库开始兴起。它们的目标是通过提供新的数据模型（参见第2章），以及内置的复制（参见第5章）和分区（参见第6章）等手段来改进传统的关系模型。然而事务却成了这场变革的受害者：很多新一代的数据库完全放弃了事务支持，或者将其重新定义，即替换为比以前弱得多的保证。

随着这种新型分布式数据库的炒作，很多人开始认为事务与可扩展性是相对立的两



面，而大规模系统为了性能与高可用性将不得不牺牲事务的支持<sup>[5,6]</sup>。但另一方面，还有一些数据库供应商则坚称事务是“关键应用”和“高价值数据”所必备的重要功能。而笔者看来，这两个观点未免都有些夸大其词。

与其他技术一样，事务有其优势，也有其自身的局限性。为了更好地理解事务设计的权衡之道，让我们考虑正常运行和各种极端（但确实存在）情况，详细分析事务可以为我们提供哪些保证。

## ACID的含义

事务所提供的安全保证即大家所熟知的ACID，分别代表原子性（Atomicity），一致性（Consistency），隔离性（Isolation）与持久性（Durability），取这四个特性的首字母。最早由Theo Härder和Andreas Reuter于1983年为精确描述数据库的容错机制而定义。

但实际上，各家数据库所实现的ACID并不尽相同。例如，我们稍后就会看到，围绕着“隔离性”就存在很多含糊不清的争议<sup>[8]</sup>。想法非常美好，细节方见真章。当听到一个系统声称自己“兼容ACID”时，其实你无法确信它究竟能提供了什么样的保证，现在的ACID更像是一个市场营销用语。

而不符合ACID标准的系统有时被冠以BASE，取另外几个特性的首字母，即基本可用性（Basically Available），软状态（Soft state）和最终一致性（Eventual consistency）<sup>[9]</sup>。听下来它似乎比ACID更加模棱两可。BASE唯一可以确定的是“它不是ACID”，此外它几乎没有承诺任何东西。

我们还是先搞清楚原子性，一致性，隔离性和持久性的准确含义，目标是建立对事务思想的清晰而牢固的认识。

### 原子性

通常，原子是指不可分解为更小粒度的东西。这个术语在计算机的不同领域里有着相似但却微妙的差异。例如，多线程编程中，如果某线程执行一个原子操作，这意味着其他线程是无法看到该操作的中间结果。它只能处于操作之前或操作之后的状态，而不是两者之间的状态。

而ACID中的原子性并不关乎多个操作的并发性，它并没有描述多个线程试图访问相同的数据会发生什么情况，后者其实是由ACID的隔离性所定义（参见本章后面的“隔离性”）。

ACID原子性其实描述了客户端发起一个包含多个写操作的请求时可能发生的情况，例如在完成了一部分写入之后，系统发生了故障，包括进程崩溃，网络中断，磁盘变满或者违反了某种完整性约束等；把多个写操作纳入到一个原子事务，万一出现了上述故障而导致没法完成最终提交时，则事务会中止，并且数据库须丢弃或撤销那些局部完成的更改。

假如没有原子性保证，当多个更新操作中间发生了错误，就需要知道哪些更改已经生效，哪些没有生效，这个寻找过程会非常麻烦。或许应用程序可以重试，但情况类似，并且可能导致重复更新或者不正确的结果。原子性则大大简化了这个问题：如果事务已经中止，应用程序可以确定没有实质发生任何更改，所以可以安全地重试。

因此ACID中原子性所定义的特征是：在出错时中止事务，并将部分完成的写入全部丢弃。也许可中止性比原子性更为准确，不过我们还是沿用原子性这个惯用术语。

## 一致性

一致性非常重要，但它在不同场景有着不同的具体含义，例如：

- 第5章我们讨论了副本一致性以及异步复制模型时，引出了最终一致性问题（参见第5章“复制滞后问题”）。
- 一致性哈希则是某些系统用于动态分区再平衡的方法（参见第6章“一致性哈希”）。
- CAP理论中，一致性一词用来表示线性化（参见第9章“可线性化”）。
- 而在ACID中，一致性主要指数据库处于应用程序所期待的“预期状态”。

可以看出，同一个词至少有四种不同的含义。

ACID中的一致性的主要是指对数据有特定的预期状态，任何数据更改必须满足这些状态约束（或者恒等条件）。例如，对于一个账单系统，账户的贷款余额应和借款余额保持平衡。如果某事务从一个有效的状态开始，并且事务中任何更新操作都没有违背约束，那么最后的结果依然符合有效状态。

这种一致性本质上要求应用层来维护状态一致（或者恒等），应用程序有责任正确地定义事务来保持一致性。这不是数据库可以保证的事情：即如果提供的数据修改违背了恒等条件，数据库很难检测进而阻止该操作（数据库可以完成针对某些特定类型的恒等约束检查，例如使用外键约束或唯一性约束。但通常主要靠应用程序来定义数据的有效/无效状态，数据库主要负责存储）。

原子性，隔离性和持久性是数据库自身的属性，而ACID中的一致性更多是应用层的属性。应用程序可能借助数据库提供的原子性和隔离性，以达到一致性，但一致性本身并不源于数据库。因此，字母C其实并不应该属于ACID<sup>注1</sup>。

### 隔离性

大多数数据库都支持多个客户端同时访问。如果读取和写入的是不同数据，这肯定没有什么问题；但如果访问相同的记录，则可能会遇到并发问题（即带来竞争条件）。

一个简单例子如图7-1所示。假设有两个客户端同时增加数据库中的一个计数器。每个客户首先读取当前值，再客户端增加1，然后写回新值（这里假设数据库尚不支持自增操作）。图7-1中，由于有两次相加，计数器应该由42增加到44，但实际上由于竞争条件最终结果却是43。

ACID语义中的隔离性意味着并发执行的多个事务相互隔离，它们不能互相交叉。经典的数据库教材把隔离定义为可串行化，这意味着可以假装它是数据库上运行的唯一事务。虽然实际上它们可能同时运行，但数据库系统要确保当事务提交时，其结果与串行执行（一个接一个执行）完全相同<sup>[10]</sup>。

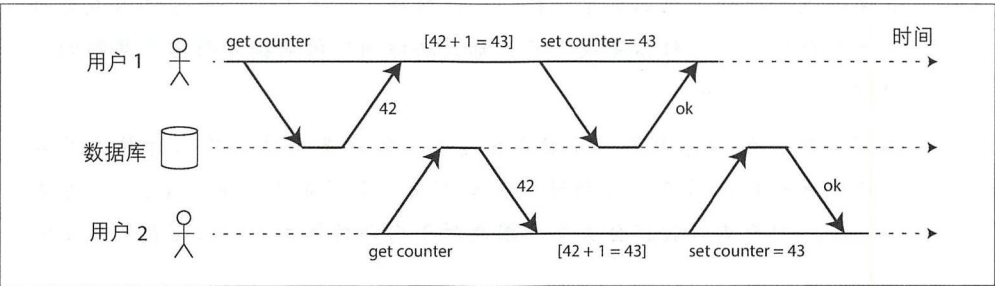


图7-1：两个客户端同时递增一个计数器产生竞争条件

然而实践中，由于性能问题很少使用串行化隔离。一些流行的数据库，如Oracle 11g，甚至根本就没有实现它。Oracle虽然也有声称“串行化”的功能，但它本质上实现的是快照隔离，后者提供了比串行化更弱的保证<sup>[8,11]</sup>。我们将在本章后面“弱隔离级别”中讨论快照隔离以及其他形式的隔离。

### 持久性

数据库系统本质上是提供一个安全可靠的地方来存储数据而不用担心数据丢失等。持

注1： Joe Hellerstein 曾经在文献<sup>[7]</sup>中评论，字母C只是为了使ACID这个缩略词听起来更为顺口，当时并非觉得这是件很重要的事情。



久性就是这样的承诺，它保证一旦事务提交成功，即使存在硬件故障或数据库崩溃，事务所写入的任何数据也不会消失。

对于单节点数据库，持久性通常意味着数据已被写入非易失性存储设备，如硬盘或SSD。在写入执行过程中，通常还涉及预写日志等（参阅第3章“可靠的B-tree”），这样万一磁盘数据损坏可以进行恢复。而对于支持远程复制的数据库，持久性则意味着数据已成功复制到多个节点。为了实现持久性的保证，数据库必须等到这些写入或复制完成之后才能报告事务成功提交。

正如第1章的“可靠性”所提到的，其实不存在完美的持久性。例如，所有的硬盘和所有的备份如果都同时被（人为）销毁了，那么数据库也无能为力。

## 复制与持久性

从历史上看，持久性最早意味着是写入磁带来存档，之后演变为写入磁盘或SSD。最近，它又代表着多节点间复制。对比起来，哪种方式更好呢？

答案是，没有任何一个可以称为完美：

- 如果写入了磁盘但机器发生死机，即使数据没有丢失，在重启机器或者将磁盘转移到另一台机器完成之前，都无法访问任何数据。而基于复制的系统则可以继续可用。
- 某些故障，例如停电或一个软件bug，可能在某个特定输入时触发了异常进而导致所有节点全部崩溃，甚至删除所有的副本（参阅第1章“可靠性”），内存中的数据会丢失，因此对于内存数据库，写入磁盘仍然是需要的。
- 在异步复制系统中，当主节点不可用时，最近的写入操作因为可能没有及时完成同步而导致更新丢失（参阅第5章“处理节点失效”）。
- 当电源突然断电时，特别对于固态硬盘可能发生意外：连fsync之后的数据也不能保证可以正确恢复<sup>[12]</sup>。和所有其他类型的软件类似，磁盘的固件也可能存在bug<sup>[13,14]</sup>。
- 考虑到数据库存储引擎与文件系统之间复杂而微妙的关系，其中也可能包含难以追踪的bug，并最终导致在系统发生崩溃后，磁盘上的文件也出现损坏<sup>[15,16]</sup>。

- 磁盘上的数据可能会悄无声息地出现损坏但又没有被及时检测到<sup>[17]</sup>。这种情况如果持续发生，则多个副本甚至最近的备份都可能已经损坏。此时需从历史备份中恢复数据。
- 一项关于固态硬盘的研究发现，在部署的前四年，30%到80%的固态硬盘至少发生一个坏块<sup>[18]</sup>。磁盘的坏道率比较低，但整盘发生完全失效的概率却比固态硬盘高。
- 某些固态硬盘如果遭遇突然断电，可能会在接下来的几周内发生丢失数据情况，温度在里面起了关键性作用<sup>[19]</sup>。

现实情况是，没有哪一项技术可以提供绝对的持久性保证。这些都是帮助降低风险的手段，应该组合使用它们，包括写入磁盘、复制到远程机器以及备份等。因此对任何理论上的“保证”一定要谨慎对待。

## 单对象与多对象事务操作

回顾一下，ACID中的原子性和隔离性主要针对客户端在同一事务中包含多个写操作时，数据库所提供的保证：

### 原子性

如果一系列写操作中间发生了错误，则事务必须中止，并且事务中已完成的写入应该被丢弃。换言之，不用担心数据库的部分失败，它总是保证要么全部成功，要么全部失败。

### 隔离性

同时运行的事务不应相互干扰。例如，如果某个事务进行多次写入，则另一个事务应该观察到的是其全部完成（或者一个都没完成）的结果，而不应该看到中间的部分结果。

这些定义假定在一个事务中会修改多个对象（如行，文档，记录等）。这种多对象事务目的通常是为了在多个数据对象之间保持同步。图7-2展示了一个电子邮件应用的例子。要显示用户的未读邮件数量，可以执行查询如下：

```
SELECT COUNT(*) FROM emails WHERE recipient_id = 2 AND unread_flag = true
```

如果电子邮件太多，你会发现查询太慢，然后决定将未读的数量直接保存在一个单独的字段中（虽然这违反了范式要求）。这样每当收到一个新邮件，需要增加未读计数器；当邮件标记为已读时，还需减少该计数器。

在图7-2中，用户2遇到些异常情况：邮箱列表已显示了未读消息，但计数器却还未更新，所显式的数目是0<sup>注2</sup>。隔离性将保证用户2看到要么是更新后的电子邮件和更新后的计数器，要么是二者都未更新，而不会是两者不一致。

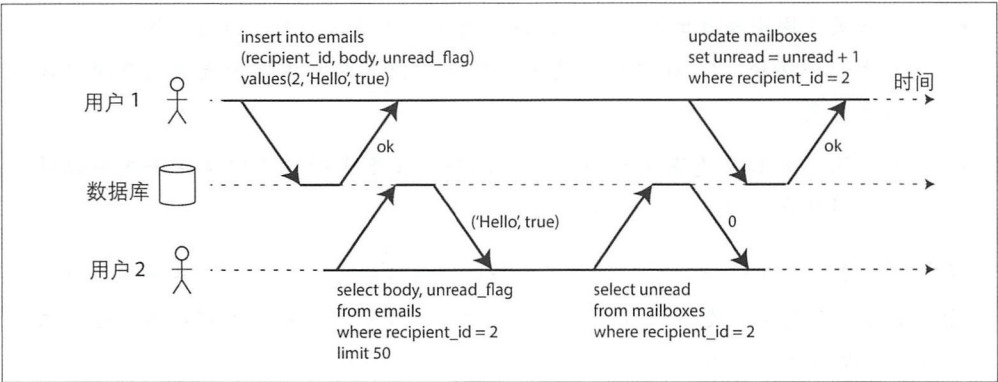


图7-2：某个事务读取了另一个事务中尚未提交的写入（“脏读”），违反了隔离性

图7-3则说明了对原子性的需求：如果事务执行过程中发生错误，导致邮箱和未读计数器二者不同步。则事务将被中止，且此之前插入的电子邮件将被回滚。

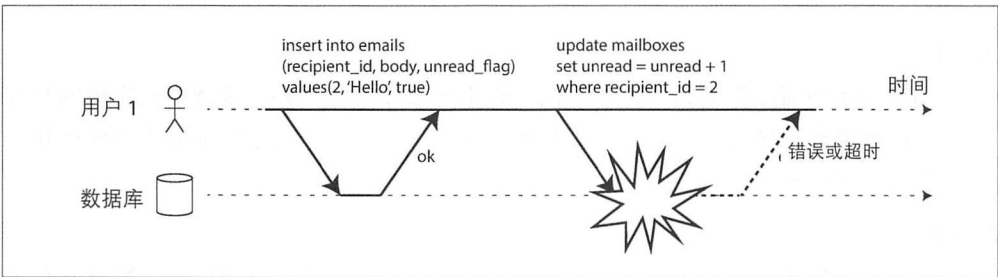


图7-3：如果发生了错误，原子性保证撤销事务中已完成的部分写入，以避免状态不一致

多对象事务要求确定知道事务包含了哪些读写操作。对于关系数据库，客户端通常与数据库服务器建立TCP网络连接，因而对于特定的某个连接，SQL语句`BEGIN TRANSACTION`和`COMMIT`之间的所有操作都属于同一个事务<sup>注3</sup>。

注2：这个电子邮件例子中所出现的计数器错误谈不上多么关键。如果需要，可以试着把未读邮件计数器替换为账户余额，把邮件替换为转账交易。

注3：这种方式并不完美。如果TCP连接中断，则事务也必须中止。现在假定中断发生在客户端提交请求之后，服务器确认提交完成之前，则客户端最后不知道该事务是否已完成提交。为了解决该问题，事务管理器可以定义一个唯一的事务标识符来逻辑上绑定一组写操作，且该事务标识符独立于与TCP链接。我们将在第12章中详细解释。



而许多非关系数据库则不会将这些操作组合在一起。即使他们可能支持多对象API（例如，键-值存储的*multi-put* API可以在一个操作中更新多个键），但并不意味着具有事务语义，例如可能出现某些键更新成功了而其他则发生了失败，最后结果是数据库处于部分更新的状态。

## 单对象写入

原子性和隔离性也同样适用于单个对象的更新。例如，假设向数据库写入20KB的JSON文档：

- 如果发送了第一个10KB之后网络连接中断，数据库是否只存储了无法完整解析的10KB JSON片段呢？
- 如果数据库在覆盖磁盘现有数据时发生电源故障，最终是否是新旧值混杂在一起？
- 如果另一个客户端在写入的过程中读取该文档，是否会看到部分更新的文档内容？

这些问题着实让人头疼，因此存储引擎几乎必备的设计就是在单节点、单个对象层面上提供原子性和隔离性（比如key-value对）。例如，出现宕机时，基于日志恢复来实现原子性（参阅第3章“可靠的B-tree”），对每个对象采用加锁的方式（每次只允许一个线程访问对象）来实现隔离。

某些数据库还提供了高级的原子操作，例如原子自增操作<sup>注4</sup>，这样就不需要像图7-1那样执行读取-修改-写回，类似地还有原子比较-设置操作，即只有当前值没有被人修改时才执行写入（参阅本章后面的“compare-and-set”）。

这些单对象操作可以有效防止多个客户端并发修改同一对象时的更新丢失问题（参阅本章后面的“防止更新丢失”）。但需要注意，它们并不是通常意义上的事务。虽然compare-and-set和其他单对象操作有时也被称为“轻量级事务”，甚至“ACID”<sup>[20-22]</sup>，但这里其实存在一些误导性。通常意义上的事务针对的是多个对象，将多个操作聚合为一个逻辑执行单元。

## 多对象事务的必要性

许多分布式数据存储系统不支持多对象事务，主要是因为当出现跨分区时，多对象事务非常难以正确实现，同时在高可用或者极致性能的场景下也会带来很多负面影响。

---

注4：严格来说，原子自增这个术语更多是指多线程编程的语义。而在ACID的背景下，称为隔离增量或者串行化增量貌似更为准确。不过，我们不打算在这里咬文嚼字。

但是，分布式数据库实现事务并非不可能，并不存什么原理上的限制，我们将在第9章讨论分布式事务的实现。

但是否所有应用都需要多对象事务呢？是否可能只用键-值数据模型和单对象操作就可以满足应用需求？

的确有一些情况，只进行单个对象的插入、更新和删除就足够了。但是，还有许多其他情况要求写入多个不同的对象并进行协调：

- 对于关系数据模型，表中的某行可能是另一个表中的外键。类似地，在图数据模型中，顶点具有多个边链接到其他的顶点。多对象事务用以确保这些外键引用的有效性，即当插入多个相互引用的记录时，保证外键总是最新、正确的，否则数据更新就变得毫无意义。
- 对于文档数据模型，如果待更新的字段都在同一个文档中，则视为单个对象，此时不需要多对象事务。但是，缺少join支持的文档数据库往往会滋生反规范化（参阅第2章“关系数据库与文档数据库现状”），如图7-2所示，当更新这种非规范化数据时，就需要一次更新多个文档。此时多对象事务就可以有效防止非规范化数据之间出现不同步。
- 对于带有二级索引的数据库（除了纯粹键-值存储以外几乎所有其他系统都支持二级索引），每次更改值时都需要同步更新索引。从事务角度来看，这些索引是不同的数据库对象：如果没有事务隔离，就会出现部分索引更新。

即使没有事务支持，或许上层应用依然可以工作，然而在没有原子性保证时，错误处理就会异常复杂，而缺乏隔离性则容易出现并发性方面的各种奇怪问题。我们将在本章后面“弱隔离级别”中讨论这些问题，并在第12章中探讨其他的一些方案。

## 处理错误与中止

事务的一个关键特性是，如果发生了意外，所有操作被中止，之后可以安全地重试。ACID数据库基于这样的理念：如果存在违反原子性、隔离性或持久性的风险，则完全放弃整个事务，而不是部分放弃。

然而，并不是所有的系统都遵循上述理念。例如，无主节点复制的数据存储（参阅第5章“无主节点复制”）会在“尽力而为”的基础上尝试多做些工作，可以概括理解为：数据库已经尽其所能，但万一遇到错误，系统并不会撤销已完成的操作，此时需要应用程序来负责从错误中进行恢复。

确实我们无法彻底避免错误，然而许多开发人员喜欢只考虑正常的处理路径，而忽视

错误处理。比如像流行的Rails ActiveRecord和Django这样的对象关系映射（ORM）框架，当事务出现异常时不会进行重试而只是简单地抛出堆栈信息，用户虽然得到了错误提示，但所有之前的输入都会被抛弃。这绝对不应该，支持安全的重试机制才是中止流程的重点。

重试中止的事务虽然是一个简单有效的错误处理机制，但它并不完美：

- 如果事务实际已经执行成功，但返回给客户端的消息在网络传输时发生意外（所以在客户端看来事务是失败），那么重试就会导致重复执行，此时需要额外的应用级重复数据删除机制。
- 如果错误是由于系统超负荷所导致，则重试事务将使情况变得更糟。为此，可以设定一个重试次数上限，例如指数回退，同时要尝试解决系统过载本身的问题。
- 由临时性故障（例如死锁，隔离违例，网络闪断和节点切换等）所导致的错误需要重试。但出个出现了永久性故障（例如违反约束），则重试毫无意义。
- 如果在数据库之外，事务还产生其他副作用，即使事务被中止，这些副作用可能已事实生效。例如，假设更新操作还附带发送一封电子邮件，肯定不希望每次重试时都发送邮件。如果想要确保多个不同的系统同时提交或者放弃，可以考虑采用两阶段提交（参阅第9章“原子提交与两阶段提交”）。
- 如果客户端进程在重试过程中也发生失败，没有其他人继续负责重试，则那些待写入的数据可能会因此而丢失。

## 弱隔离级别

如果两个事务操作的是不同的数据，即不存在数据依赖关系，则它们可以安全地并行执行。只有出现某个事务修改数据而另一个事务同时要读取该数据，或者两个事务同时修改相同数据时，才会引发并发问题（引入了竞争条件）。

并发性相关的错误很难通过测试发现，这类错误通常只在某些特定时刻才会触发，这种时机相关的问题发生概率低，稳定重现比较困难。并发性也很难对其进行推理分析，特别是对于一个大型应用程序，几乎不可能知道哪些代码正在访问数据库；只有一个用户访问数据时，程序开发就足够困难了，当出现多用户并发时情况会变得更加复杂，每一块数据随时都可能被多个用户所修改。

正因如此，数据库一直试图通过事务隔离来对应用开发者隐藏内部的各种并发问题。从理论上讲，隔离是假装没有发生并发，让程序员的生活更轻松，而可串行化隔离意



意味着数据库保证事务的最终执行结果与串行（即一次一个，没有任何并发）执行结果相同。

实现隔离绝不是想象的那么简单。可串行化的隔离会严重影响性能，而许多数据库却不愿意牺牲性能<sup>[8]</sup>，因而更多倾向于采用较弱的隔离级别，它可以防止某些但并非全部的并发问题。这些弱隔离级别理解起来更为困难，甚至可能会带来一些难以捉摸的隐患，但在实践中还是被广泛使用<sup>[23]</sup>。

弱隔离所引发的并发性错误绝非仅是理论存在，它们已经造成了大量的资金损失<sup>[24,25]</sup>，审计部门的调查<sup>[26]</sup>，以及客户数据破坏<sup>[27]</sup>等。对此，有一个流行的说法是“如果你正在处理财务数据，请上ACID系统！”，但这样的建议其实没有太大实际意义，因为很多流行的关系数据库系统（通常被认为是“ACID兼容”）其实也采用的是弱级别隔离，所以它们未必可以阻止类似错误的发生。

与其盲目地相信这些宣传，不如对存在的并发问题以及如何防范有一个全面、深刻的理解。然后，我们就可以使用所掌握的工具和方法来构建正确、可靠的应用。

本节将分析几个实际中经常用到的弱级别（非串行化）隔离，并详细讨论可能（或者不可能）发生的竞争条件，有了这些认识之后，可以帮助判断自己的应用更适合什么样的隔离级别。在下一节，我们将详细讨论可串行化。我们主要以示例的方式讨论隔离级别，如果你需要的是严格、正式的定义和分析，可以参考文献[28-30]。

## 读-提交

读-提交是最基本<sup>注5</sup>的事务隔离级别，它只提供以下两个保证：

1. 读数据库时，只能看到已成功提交的数据（防止“脏读”）。
2. 写数据库时，只会覆盖已成功提交的数据（防止“脏写”）。

这两个保证更深入的介绍如下：

### 防止脏读

假定某个事务已经完成部分数据写入，但事务尚未提交（或中止），此时另一个事务是否可以看到尚未提交的数据呢？如果是的话，那就是脏读<sup>[2]</sup>。

读-提交级别的事务隔离必须做到防止发生脏读。这意味着事务的任何写入只有在成

---

注5：某些数据库甚至提供更弱的隔离级别，称为读-未提交。它只防止脏写，而不防止脏读。

功提交之后，才能被其他人观察到（并且所有的写全部可见）。如图7-4所示，用户1设置了 $x = 3$ ，在用户1的事务未提交之前，用户2的`get x`操作依旧返回的是旧值2。

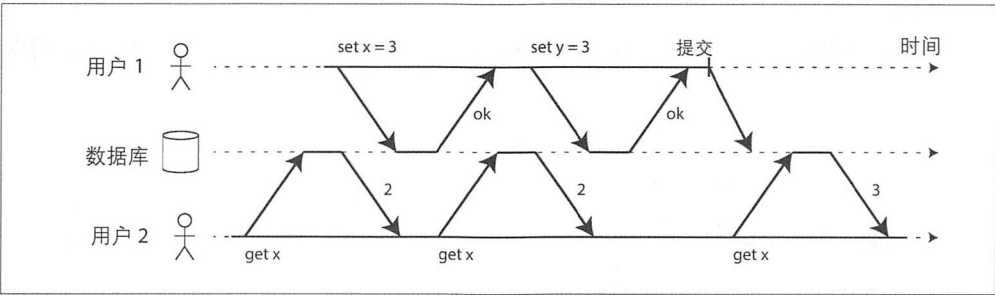


图7-4：没有脏读时，用户2只有在用户1的事务提交之后才能看到 $x$ 的新值

当有以下需求时，需要防止脏读：

- 如果事务需要更新多个对象，脏读意味着另一个事务可能会看到部分更新，而非全部。例如，图7-2就是一个电子邮件应用的脏读例子，用户看到了新的未读电子邮件，但看不到更新的计数器。观察到部分更新的数据可能会造成用户的困惑，并由此引发一些不必要的后续操作。
- 如果事务发生中止，则所有写入操作都需要回滚（见图7-3）。如果发生了脏读，这意味着它可能会看到一些稍后被回滚的数据，而这些数据并未实际提交到数据库中。之后所引发的后果可能都会变得难以预测。

### 防止脏写

如果两个事务同时尝试更新相同的对象，会发生什么情况呢？我们不清楚写入的顺序，但可以想象后写的操作会覆盖较早的写入。

但是，如果先前的写入是尚未提交事务的一部分，是否还是被覆盖？如果是，那就是脏写<sup>[28]</sup>。读-提交隔离级别下所提交的事务可以防止脏写，通常的方式是推迟第二个写请求，直到前面的事务完成提交（或者中止）。

防止脏写可以避免以下并发问题：

- 如果事务需要更新多个对象，脏写会带来非预期的错误结果。例如，考虑图7-5的二手车销售网站，Alice和Bob两个人试图购买同一辆车。而购买汽车需要两次数据库写入：网站上商品买主需要更新为新买家，销售发票也需要随之更新。对于图7-5的例子，车主被改为Bob（因为他成功地抢先更新了车辆表单），但发

票却发给了Alice（因为她成功的先执行了发票表单）。读提交隔离要防止这种事故。

- 但是，读-提交隔离并不能解决图7-1中计数器增量的竞争情况。对于后者，第二次写入确实在第一个事务提交后才执行，虽然不属于脏写，但结果仍然是错误的。在接下来的“防止更新丢失”中，我们将讨论如何安全递增计数器。

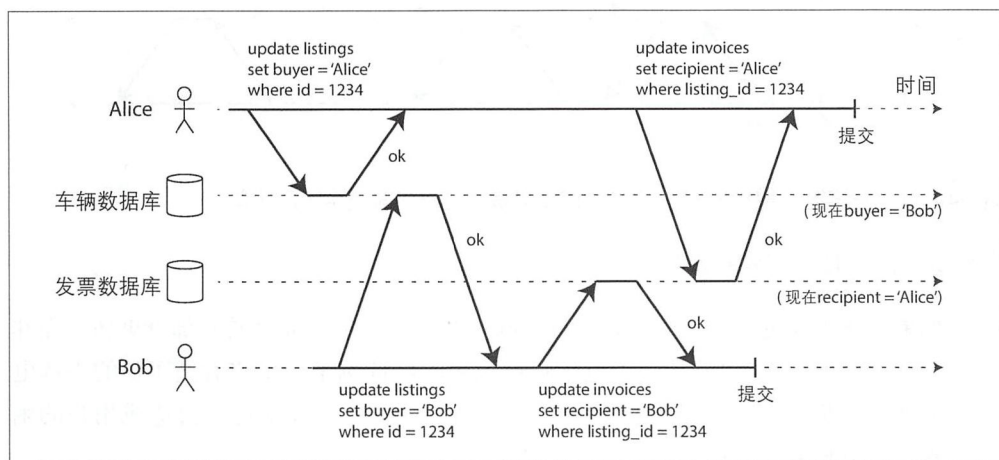


图7-5：脏写问题导致不同事务的并发写入最终混杂在一起

## 实现读-提交

读-提交隔离非常流行，它是Oracle 11g、PostgreSQL、SQL Server 2012、MemSQL，以及许多其他数据库的默认配置<sup>[8]</sup>。

数据库通常采用行级锁来防止脏写：当事务想修改某个对象（例如行或文档）时，它必须首先获得该对象的锁；然后一直持有锁直到事务提交（或中止）。给定时刻，只有一个事务可以拿到特定对象的锁，如果有另一个事务尝试更新同一个对象，则必须等待，直到前面的事务完成了提交（或中止）后，才能获得锁并继续。这种锁定是由处于读-提交模式（或更强的隔离级别）数据库自动完成的。

那如何防止脏读呢？一种选择是使用相同的锁，所有试图读取该对象的事务必须先申请锁，事务完成后释放锁。从而确保不会发生读取一个脏的、未提交的值（因为锁在那段期间一直由一个事务所持有）。

然而，读锁的方式在实际中并不可行，因为运行时间较长的写事务会导致许多只读的事务等待太长时间，这会严重影响只读事务的响应延迟，且可操作性差：由于读锁，应用程序任何局部的性能问题会扩散进而影响整个应用，产生连锁反应。



因此，大多数数据库<sup>注6</sup>采用了图7-4所示的方法来防止脏读：对于每个待更新的对象，数据库都会维护其旧值和当前持锁事务将要设置的新值两个版本。在事务提交之前，所有其他读操作都读取旧值；仅当写事务提交之后，才会切换到读取新值。

## 快照级别隔离与可重复读

表面上看读-提交级别隔离，可能会认为它已经满足了事务所需要一切特征：它支持中止（原子性所必须的），可以防止读取不完整的结果，并且防止并发写的混合。事实上，这些确实非常有用，相比没有事务的系统，它的确提供了更多的保证。

但是，在使用此隔离级别时，仍然有很多场景可能导致并发错误。如图7-6所示。

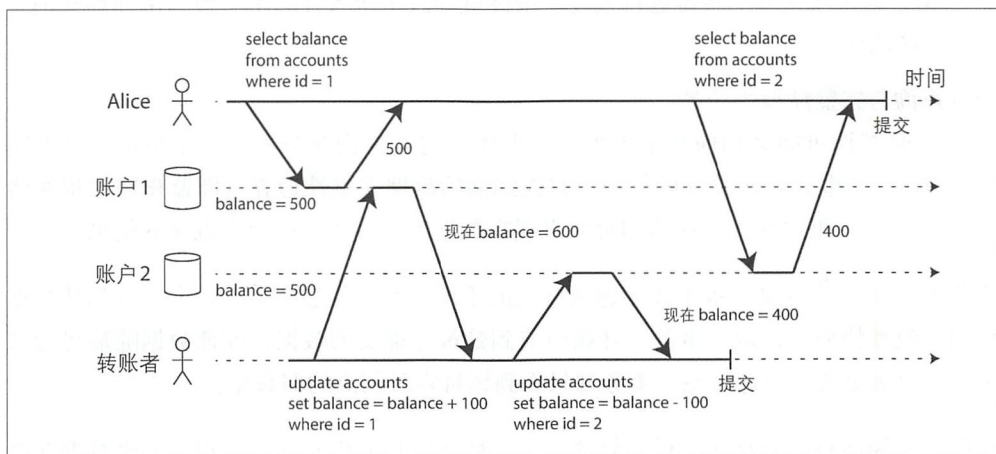


图7-6: 读倾斜: Alice观察数据库处于不一致的状态

假设Alice在银行有1000美元的存款，分为两个账户，每个500美元。现在有这样一笔转账交易从账户1转100美元到账户2。如果在她提交转账请求之后而银行数据库系统执行转账的过程中，来查看两个账户的余额，她有可能会看到账号2在收到转账之前的余额（500美元），和账户1在完成转账之后的余额（400美元）。对于Alice来说，貌似她的账户总共只有900美元，有100美元消失了。

这种异常现象被称为不可重复读取（nonrepeatable read）或读倾斜（read skew）。如果Alice在交易结束时再次读取账户1的余额，她将看到不同的值（600美元）。读倾斜在读-提交隔离语义下是可以接受的，Alice所看到的账户余额的确都是账户当时的最新值。

注6: 编写本书时,唯一还使用锁的主流数据库是设置`read_committed_snapshot = off`时的IBM DB2和Microsoft SQL Server<sup>[23,36]</sup>。



倾斜 (skew) 这个词有些滥用了, 我们以前使用它是由于热点而导致负载不平衡 (参阅第6章 “倾斜的负载与减轻热点”), 而这里则主要是指时间异常。

对于Alice这个例子, 这并非一个永久性问题, 例如几秒钟之后当她重新加载银行页面, 可能就能看到一致的账户余额。但是, 还有些场景则不能容忍这种暂时的不一致:

### 备份场景

备份任务要复制整个数据库, 这可能需要数小时才能完成。在备份过程中, 可以继续写入数据库。因此, 得到镜像里可能包含部分旧版本数据和部分新版本数据。如果从这样的备份进行恢复, 最终就导致了永久性的不一致 (例如那些消失的存款)。

### 分析查询与完整性检查场景

有时查询可能会扫描几乎大半个数据库。这类查询在分析业务中很常见 (参阅第3章 “事务处理或分析”), 亦或定期的数据完整性检查 (即监视数据损坏情况)。如果这些查询在不同时间点观察数据库, 可能会返回无意义的结果。

快照级别隔离<sup>[28]</sup>这是阶级上述问题最常见的手段。其总体想法是, 每个事务都从数据库的一致性快照中读取, 事务一开始所看到是最近提交的数据, 即使数据随后可能被另一个事务更改, 但保证每个事务都只看到该特定时间点的旧数据。

快照级别隔离对于长时间运行的只读查询 (如备份和分析) 非常有用。如果数据在执行查询的同时还在发生变化, 那么查询结果对应的物理含义就难以理清。而如果查询的是数据库在某时刻点所冻结的一致性快照, 则查询结果的含义非常明确。

快照级别隔离非常流行, 目前PostgreSQL, MySQL的InnoDB存储引擎, Oracle, SQL Server等<sup>[23,31,32]</sup>都已经支持。

## 实现快照级别隔离

与读-提交隔离类似, 快照级别隔离的实现通常采用写锁来防止脏写 (参阅本章前面的 “实现读-提交”), 这意味着正在进行写操作的事务会阻止同一对象上的其他事务。但是, 读取则不需要加锁。从性能角度看, 快照级别隔离的一个关键点是读操作不会阻止写操作, 反之亦然。这使得数据库可以在处理正常写入的同时, 在一致性快照上执行长时间的只读查询, 且两者之间没有任何锁的竞争。

为了实现快照级别隔离, 数据库采用了一种类似于图7-4中防止脏读但却更为通用的

机制。考虑到多个正在进行的事务可能会在不同的时间点查看数据库状态，所以数据库保留了对象多个不同的提交版本，这种技术因此也被称为多版本并发控制（Multi-Version Concurrency Control，MVCC）。

如果只是为了提供读-提交级别隔离，而不是完整的快照级别隔离，则只保留对象的两个版本就足够了：一个已提交的旧版本和尚未提交的新版本。所以，支持快照级别隔离的存储引擎往往直接采用MVCC来实现读-提交隔离。典型的做法是，在读-提交级别下，对每一个不同的查询单独创建一个快照；而快照级别隔离则是使用一个快照来运行整个事务。

图7-7说明了PostgreSQL如何实现基于MVCC的快照级别隔离<sup>[31]</sup>（其他实现基本类似）。当事务开始时，首先赋予一个唯一的、单调递增的事务ID<sup>注7</sup>（txid）。每当事务向数据库写入新内容时，所写的数据库都会被标记写入者的事务ID。

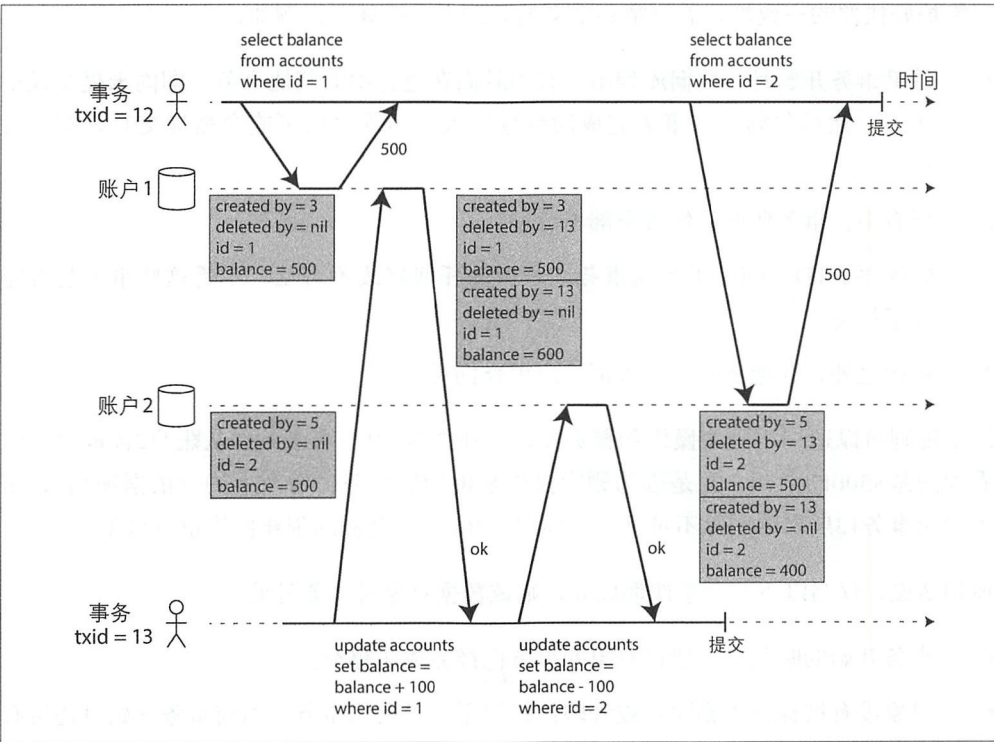


图7-7：采用多版本技术实现快照级别隔离

注7：事务ID其实是个32位整数，大约在40亿次事务之后会溢出。PostgreSQL的vacuum进程会执行后台清理，确保不会出现异常情况。



表中的每一行都有一个`created_by`字段，其中包含了创建该行的事务ID。每一行还有一个`deleted_by`字段，初始为空。如果事务要删除某行，该行实际上并未从数据库中删除，而只是将`deleted_by`字段设置为请求删除的事务ID（仅仅标记为删除）。事后，当确定没有其他事务引用该标记删除的行时，数据库的垃圾回收进程才去真正删除并释放存储空间。

这样一笔更新操作在内部会被转换为一个删除操作加一个创建操作。例如，图7-7中，事务13从账户2中扣除100美元，余额从500美元减为400美元。`accounts`表里会出现两行账户2：一个余额为\$500但标记为删除的行（由事务13删除），另一个余额为\$400，由事务13创建。

## 一致性快照的可见性规则

当事务读数据库时，通过事务ID可以决定哪些对象可见，哪些不可见。要想对上层应用维护好快照的一致性，需要精心定义数据的可见性规则。例如：

1. 每笔事务开始时，数据库列出所有当时尚在进行中的其他事务（即尚未提交或中止），然后忽略这些事务完成的部分写入（尽管之后可能会被提交），即不可见。
2. 所有中止事务所做的修改全部不可见。
3. 较晚事务ID（即晚于当前事务）所做的任何修改不可见，不管这些事务是否完成了提交。
4. 除此之外，其他所有的写入都对应用查询可见。

以上规则可以适用于创建操作和删除操作。在图7-7中，当事务12从账户2读取时，它看到的是\$500的余额，这是因为删除操作是由稍后事务13所产生的（依据规则3，事务12对事务13所做的删除不可见），同理，400美元余额的创建操作也不可见。

换句话说，仅当以下两个条件都成立，则该数据对象对事务可见：

- 事务开始的时刻，创建该对象的事务已经完成了提交。
- 对象没有被标记为删除；或者即使标记了，但删除事务在当前事务开始时还没有完成提交。

长时间运行的事务可能会使用快照很长时间，从其他事务的角度来看，它可能在持续访问正在被覆盖或删除的内容。由于没有就地更新，而是每次修改总创建一个新版本，因此数据库可以以较小的运行代价来维护一致性快照。

## 索引与快照级别隔离

接下来一个问题是，这种多版本数据库该如何支持索引呢？一种方案是索引直接指向对象的所有版本，然后想办法过滤对当前事务不可见的那些版本。当后台的垃圾回收进程决定删除某个旧对象版本时，对应的索引条目也需要随之删除。

在实践中，有许多细节决定了多版本并发控制的实际性能表现。例如，可以把同一对象的不同版本放在一个内存页面上，PostgreSQL采取这样的优化措施来避免更新索引<sup>[31]</sup>。

CouchDB、Datomic和LMDB则使用另一种方法。它们主体结构是B-tree（参阅第3章“B-tree”），但采用了一种追加/写时复制的技术，当需要更新时，不会修改现有的页面，而总是创建一个新的修改副本，拷贝必要的内容，然后让父结点，或者递归向上直到树的root 结点都指向新创建的结点。那些不受更新影响的页面都不需要复制，保持不变并被父结点所指向<sup>[33-35]</sup>。

这种采用追加式的B-tree，每个写入事务（或一批事务）都会创建一个新的B-tree root，代表该时刻数据库的一致性快照。这时就没有必要根据事务ID再去过滤掉某些对象，每笔写入都会修改现有的B-tree，因为之后的查询可以直接作用于特定快照B-tree（有利于查询性能）。采用这种方法依然需要后台进程来执行压缩和垃圾回收。

## 可重复读与命名混淆

快照级别隔离对于只读事务特别有效。但是，具体到实现，许多数据库却对它有着不同的命名。Oracle称之为可串行化，PostgreSQL和MySQL则称为可重复读<sup>[23]</sup>。

这种命名混淆的原因是SQL标准并没有定义快照级别隔离，而仍然是基于老的System R 1975年所定义的隔离级别<sup>[2]</sup>，而当时还没有出现快照级别隔离。标准定义的是“可重复读”，这看起来比较接近于快照级别隔离，所以 PostgreSQL和MySQL称它们的快照级别隔为“可重复读”，这符合标准要求（即合规性）。

然而必须指出，SQL标准对隔离级别的定义还是存在一些缺陷，某些定义模棱两可，不够精确，且不能做到与实现无关<sup>[28]</sup>。尽管有几个数据库实现了可重复读，表面上看符合标准，但它们实际所提供的保证却大相径庭<sup>[23]</sup>。可重复读有一个更为严谨的定义，然而大多数实现并没有遵循它<sup>[29,30]</sup>。最后还要指出，IBM DB2的“可重复读”实则是可串行化级别隔离<sup>[8]</sup>。

现在的结果是，我们已经搞不清楚“可重复读”究竟代表什么了。

## 防止更新丢失

总结一下，我们所讨论的读-提交和快照级别隔离主要都是为了解决只读事务遇到并发写时可以看到什么（虽然中间也涉及脏写问题），总体而言我们还没有触及另一种情况，即两个写事务并发，而脏写只是写并发的一个特例。

写事务并发还会带来其他一些值得关注的冲突问题，最著名的就是更新丢失问题，前面图7-1正是这样的例子。

更新丢失可能发生在这样一个操作场景中：应用程序从数据库读取某些值，根据应用逻辑做出修改，然后写回新值（*read-modify-write*过程）。当有两个事务在同样的数据对象上执行类似操作时，由于隔离性，第二个写操作并不包括第一个事务修改后的值，最终会导致第一个事务的修改值可能会丢失。这种冲突还可能在其他不同的场景下发生，例如：

- 递增计数器，或更新账户余额（需要读取当前值，计算新值并写回更新后的值）。
- 对某复杂对象的一部分内容执行修改，例如对JSON文档中一个列表添加新元素（需要读取并解析文档，执行更改并写回修改后的文档）。
- 两个用户同时编辑wiki页面，且每个用户都尝试将整个页面发送到服务器，覆盖数据库中现有内容以使更改生效。

并发写事务冲突是一个普遍问题，目前有多种可行的解决方案。

### 原子写操作

许多数据库提供了原子更新操作，以避免在应用层代码完成“读-修改-写回”操作，如果支持的话，通常这就是最好的解决方案。例如，以下指令在多数关系数据库中都是并发安全的：

```
UPDATE counters SET value = value + 1 WHERE key = 'foo';
```

类似地，像MongoDB这样的文档数据库支持对JSON文档的某部分进行本地修改的原子操作，Redis也提供了对特定数据结构（如优先级队列）修改的原子操作。然而，并非所有的应用更新操作都可以以原子操作的方式来表达，例如维基页面的更新涉及各种文本编辑<sup>注8</sup>。无论如何，如果原子操作可行，那么它就是推荐的最佳方式。

---

注8： 尽管相当复杂，但将文本文档的编辑表示为一系列的原子操作事件流也是可行的。请参阅第5章“自动冲突解决”。



原子操作通常采用对读取对象加独占锁的方式来实现，这样在更新被提交之前不会其他事务可以读它。这种技术有时被称为游标稳定性<sup>[36,37]</sup>。另一种实现方式是强制所有的原子操作都在单线程上执行。

不过，基于对象关系映射（ORM）框架可以很容易地就产生出来非“读-修改-写回”的应用层代码，导致无法使用数据库所提供的原子操作<sup>[38]</sup>。假如你清楚知道自己在做什么，或许这并不会引发什么问题，但往往这种情况会埋下很多难以发现的潜在错误。

## 显式加锁

如果数据库不支持内置原子操作，另一种防止更新丢失的方法是由应用程序显式锁定待更新的对象。然后，应用程序可以执行“读-修改-写回”这样的操作序列；此时如果有其他事务尝试同时读取对象，则必须等待当前正在执行的序列全部完成。

例如，考虑一个多人游戏，其中几个玩家可以同时移动同一个数字。只靠原子操作可能还不够，因为应用程序还需要确保玩家的移动还需遵守其他游戏规则，这涉及一些应用层逻辑，不可能将其剥离转移给数据库层在查询时执行。此时，可以使用锁来防止两名玩家同时操作相同的棋子，参见示例7-1。

示例7-1：显式锁定行以防丢失更新

```
BEGIN TRANSACTION;

SELECT * FROM figures
  WHERE name = 'robot' AND game_id = 222
  FOR UPDATE; ❶

-- Check whether move is valid, then update the position
-- of the piece that was returned by the previous SELECT.
UPDATE figures SET position = 'c4' WHERE id = 1234;

COMMIT;
```

❶ FOR UPDATE指令指示数据库对返回的所有结果行要加锁。

首先该方法是可行的，但要做到这一点，需要仔细考虑清楚应用层的逻辑。很多代码会忘记在必要的地方加锁，结果很容易引入竞争冲突。

## 自动检测更新丢失

原子操作和锁都是通过强制“读-修改-写回”操作序列串行执行来防止丢失更新。另一种思路则是先让他们并发执行，但如果事务管理器检测到了更新丢失风险，则会中止当前事务，并强制回退到安全的“读-修改-写回”方式。



该方法的一个优点是数据库完全可以借助快照级别隔离来高效地执行检查。的确，PostgreSQL的可重复读，Oracle的可串行化以及SQL Server的快照级别隔离等，都可以自动检测何时发生了更新丢失，然后会中止违规的那个事务。但是，MySQL/InnoDB的可重复读却并不支持检测更新丢失<sup>[23]</sup>。有一些观点认为<sup>[28,30]</sup>，数据库必须防止更新丢失，要不然就不能宣称符合快照级别隔离，如果基于这样的定义，那么MySQL就属于没有完全支持快照级别隔离。

更新丢失检测是一个非常赞的功能，应用层代码因此不用依赖于某些特殊的数据库功能。开发者可能会不小心忘记使用锁或原子操作，但更新丢失检测会自动生效，有效地避免这类错误。

## 原子比较和设置

在不提供事务支持的数据库中，有时你会发现它们支持原子“比较和设置”操作（之前“单对象写入”有提到）。使用该操作可以避免更新丢失，即只有在上次读取的数据没有发生变化时才允许更新；如果已经发生了变化，则回退到“读-修改-写回”方式。

例如，为了防止两个用户同时更新同一个wiki页面，可以尝试下面的操作，这样只有当页面从上次读取之后没发生变化时，才会执行当前的更新：

```
-- This may or may not be safe, depending on the database implementation
UPDATE wiki_pages SET content = 'new content'
WHERE id = 1234 AND content = 'old content';
```

如果内容已经有了变化且值与“旧内容”不匹配，则更新将失败，需要应用层再次检查并在必要时进行重试。需要注意，如果WHERE语句是运行在数据库的某个旧的快照上，即使另一个并发写入正在运行，条件可能仍然为真，最终可能无法防止更新丢失问题。所以在使用之前，应该首先仔细检查“比较-设置”操作的安全运行条件。

## 冲突解决与复制

对于支持多副本的数据库（参见第5章），防止丢失更新还需要考虑另一个维度：由于多节点上的数据副本，不同的节点可能会并发修改数据，因此必须采取一些额外的措施来防止丢失更新。

加锁和原子修改都有个前提即只有一个最新的数据副本。然而，对于多主节点或者无主节点的多副本数据库，由于支持多个并发写，且通常以异步方式来同步更新，所以会出现多个最新的数据副本。此时加锁和原子比较将不再适用（我们将在第9章“线性化”详细讨论这个问题）。



正如第5章“检测并发写”所描述的，多副本数据库通常支持多个并发写，然后保留多个冲突版本（互称为兄弟），之后由应用层逻辑或依靠特定的数据结构来解决、合并多版本。

如果操作可交换（顺序无关，在不同的副本上以不同的顺序执行时仍然得到相同的结果），则原子操作在多副本情况下也可以工作。例如，计数器递增或向集合中添加元素等都是典型的可交换操作。这也是Riak 2.0新数据类型的设计思路，当一个值被不同的客户端同时更新时，Riak自动将更新合并在一起，避免发生更新丢失<sup>[39]</sup>。

而“最后写入获胜（LWW）”（详见第5章）冲突解决方法则容易丢失更新。不幸的是，目前LWW是许多多副本数据库的默认配置。

## 写倾斜与幻读

当多个事务同时写入同一对象时引发了两种竞争条件，即前面章节所讨论的脏写和更新丢失。为了避免数据不一致，需要借助数据库的一些内置机制，或者采取手动加锁、执行原子操作等。

然而，这还不是并发写所引发的全部问题。本节马上将看到更为微妙的写冲突的例子。

首先，设想这样一个例子：你正在开发一个应用程序来帮助医生管理医院的轮班。通常，医院会安排多个医生值班，医生也可以申请调整班次（例如他们自己生病了），但前提是确保至少一位医生还在该班次中值班<sup>[40,41]</sup>。

现在情况是，Alice和Bob是两位值班医生。两人碰巧都感到身体不适，因而都决定请假。不幸的是，他们几乎同一时刻点击了调班按钮。接下来发生的事情如图7-8所示。

每笔事务总是首先检查是否至少有两名医生目前在值班。如果是的话，则有一名医生可以安全里离开。由于数据库正在使用快照级别隔离，两个检查都返回有两名医生，所以两个事务都安全地进入到下一个阶段。接下来Alice更新自己的值班记录为离开，同样，Bob也更新自己的记录。两个事务都成功提交，最后的结果却是没有任何医生在值班，显然这违背了至少一名医生值班的业务要求。

### 定义写倾斜

这种异常情况称为写倾斜<sup>[28]</sup>。它既不是一种脏写，也不是更新丢失，两笔事务更新的是两个不同的对象（分别是Alice和Bob的值班记录）。这里的写冲突并不那么直接，





但很显然这的确是某种竞争状态：试想，如果两笔事务是串行执行，则第二个医生的申请肯定被拒绝；只有同时执行两个事务时才会引发该异常。

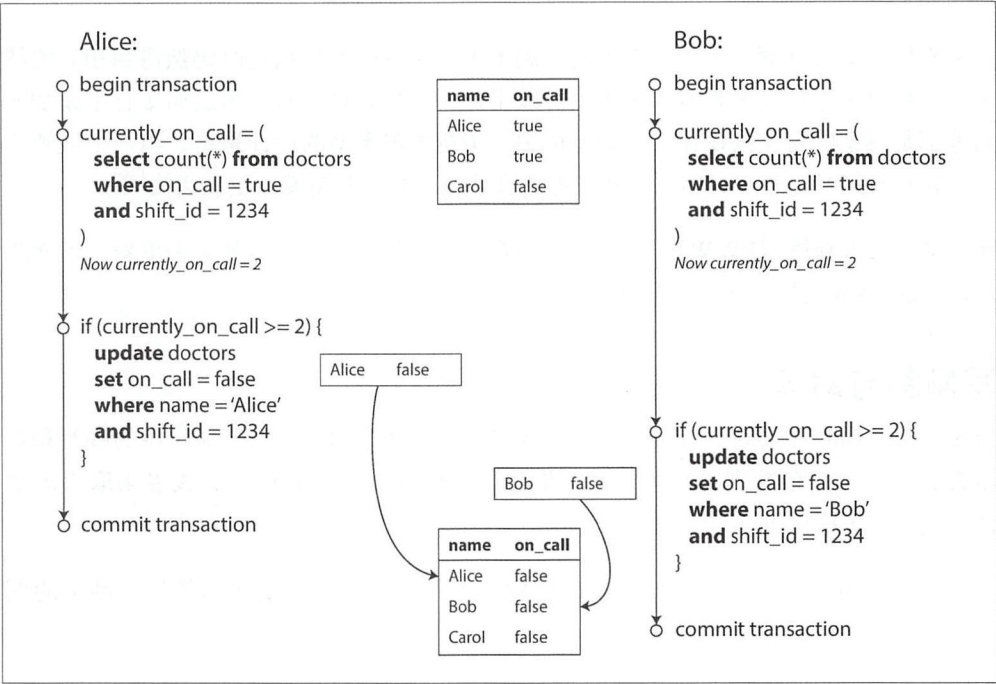


图7-8：由于写倾斜导致应用程序出现错误的示例

可以将写倾斜视为一种更广义的更新丢失问题。即如果两个事务读取相同的一组对象，然后更新其中一部分：不同的事务可能更新不同的对象，则可能发生写倾斜；而不同的事务如果更新的是同一个对象，则可能发生脏写或更新丢失（具体取决于时间窗口）。

我们已经给出了多种防范更新丢失的手段。然而对于写倾斜，可选的方案则有很多限制：

- 由于涉及多个对象，单对象的原子操作不起作用。
- 基于快照级别隔离来实现更新丢失自动检测也有问题：目前所有的数据库实现，包括PostgreSQL的可重复读，MySQL/InnoDB可重复读，Oracle可串行化以及SQL Server的快照级别隔离级别<sup>[23]</sup>都不支持检测写倾斜问题。自动防止写倾斜要求真正的可串行化隔离（参阅本章后面“可串行化”）。
- 某些数据库支持自定义约束条件，然后由数据库代为检查、执行约束（例如，唯

一性，外键约束或限制一些特定值）。但是，至少一名医生值班这样的要求涉及对多个对象进行约束，目前大多数数据库不支持这种类型约束，所以取决于具体的数据库，开发者可能可以采用触发器或物化视图来自自己实现类似约束<sup>[42]</sup>。

- 如果不能使用可串行化级别隔离，一个次优的选择是对事务依赖的行来显式的加锁。对于上述医生值班的例子，可以这样：

```
BEGIN TRANSACTION;

SELECT * FROM doctors
  WHERE on_call = true
  AND shift_id = 1234 FOR UPDATE; ❶

UPDATE doctors
  SET on_call = false
  WHERE name = 'Alice'
  AND shift_id = 1234;

COMMIT;
```

- ❶ “FOR UPDATE”语句会通知数据库对返回的所有结果行自动加锁。

## 更多写倾斜的例子

写倾斜可能看起来很晦涩拗口，可一旦深刻意识到问题的本质，就会注意到还有更多可能发生的场景。下面就是一些例子：

### 会议室预订系统

假设要求同一时间、同一个会议室不能被预订两次<sup>[43]</sup>。当有人想要预订时，首先检查是否有冲突的预订（即对同一房间的预订存在时间范围重叠），如果没有，则提交申请，参见示例7-2<sup>注9</sup>。

示例7-2：会议室预订系统，试图避免重复预订（但在快照级别隔离下不安全）

```
BEGIN TRANSACTION;

-- Check for any existing bookings that overlap with the period of noon-1pm
SELECT COUNT(*) FROM bookings
  WHERE room_id = 123 AND
    end_time > '2015-01-01 12:00' AND start_time < '2015-01-01 13:00';

-- If the previous query returned zero:
INSERT INTO bookings
  (room_id, start_time, end_time, user_id)
  VALUES (123, '2015-01-01 12:00', '2015-01-01 13:00', 666);

COMMIT;
```

注9：对于PostgreSQL，可以使用范围类型更优雅地完成此任务，但其他数据库并没有此类支持。



需要指出，快照级别隔离无法阻止并发用户预订同一个会议室。为了保证预订不会产生冲突，需要可串行化的隔离。

### 多人游戏

在示例7-1中，我们使用加锁来防止更新丢失（即两个玩家不能同时移动同一个数字）。但是，锁并不能阻止玩家将两个不同的数字移动到棋盘上的同一个位置上，或者其他可能违反游戏规则的移动。这取决于具体的游戏规则，可能需要更多的条件约束，否则很容易发生写倾斜。

### 声明一个用户名

网站通常要求每个用户有唯一的用户名，两个用户可能同时尝试创建相同的用户名。可以采用事务的方式首先来检查名称是否被使用，如果没有，则使用该名称创建账户。但是，和之前的例子类似，在快照级别隔离下这是不安全的。不过，对于该例子，一个简单的方案是采用唯一性约束（第二个事务由于违反约束，会中止创建相同用户名）。

### 防止双重开支

支付或积分相关的服务通常需要检查用户的花费不能超过其限额。一种方式是在用户的账户中插入一个临时的支出项目，然后对于所有项目，检查总开销并对比限额<sup>[44]</sup>。由于写倾斜问题，可能会同时插入两个支出项目，两个交易各自都不超额，也不会注意到对方，但加在一起则会超额。

## 为何产生写倾斜

上述所有写倾斜的例子都遵循以下类似的模式：

1. 首先输入一些匹配条件，即采用SELECT查询所有满足条件的行（例如，至少有两名医生正在值班，同一时刻房间没有预订，棋盘的某位置没有出现数字，用户名还没有被占用，账户里还有余额等）。
2. 根据查询的结果，应用层代码来决定下一步的操作（有可能继续，或者报告错误并中止）。
3. 如果应用程序决定继续执行，它将发起数据库写入（INSERT，UPDATE或DELETE）并提交事务。

而这个写操作会改变步骤2做出决定的前提条件。换句话说，如果提交写入之后再重复执行步骤1的SELECT查询，就会返回完全不同的结果，原因是刚刚的写操作改变了决定的前提条件（现在只有一名医生在值班，现在会议室已被预订，现在棋盘位置已经出现了数字，现在用户名已被占用，现在余额已经不足等）。





上述步骤可能有不同的执行顺序，例如，可以先写入，然后是SELECT查询，最后根据查询来决定是否提交或者放弃。

对于医生值班的例子，步骤3中所修改的行恰好是步骤1查询结果的一部分，因此如果先修改值班记录并加锁（SELECT FOR UPDATE），再查询可以保证事务安全，避免写倾斜。然而，对于其他例子则并不适用，它们检查的是不满足给定搜索条件的行（预期结果为空），接下来添加符合条件的行。如果步骤1的查询根本没有返回任何行，则SELECT FOR UPDATE也就无从加锁。

这种在一个事务中的写入改变了另一个事务查询结果的现象，称为幻读<sup>[3]</sup>。快照级别隔离可以避免只读查询时的幻读，但是对于我们上面所讨论那些读-写事务，它却无法解决棘手的写倾斜问题。

## 实体化冲突

如果问题的关键是查询结果中没有对象（空）可以加锁，或许可以人为引入一些可加锁的对象？

例如，对于会议室预订的例子，构造一个时间-房间表，表的每一行对应于特定时间段（例如最小15分钟间隔）的特定房间。我们提前，例如对接下来的6个月，创建好所有可能的房间与时间的组合。

现在，预订事务可以查询并锁定（SELECT FOR UPDATE）表中与查询房间和时间段所对应的行。加锁之后，即可检查是否有重叠，然后像之前一样插入新的预订。注意，这种附加表格并不存储预订相关的信息，它仅仅用于方便加锁，防止同一房间和时间范围内的重复预订。

这种方法称为实体化冲突（或物化冲突），它把幻读问题转变为针对数据库中一组具体行的锁冲突问题<sup>[11]</sup>。然而，弄清楚如何实现实体化往往也具有挑战性，实现过程也容易出错，这种把一个并发控制机制降级为数据模型的思路总是不够优雅。出于这些原因，除非万不得已，没有其他可选方案，我们不推荐采用实体化冲突。而在大多数情况下，可串行化隔离方案更为可行。

## 串行化

我们已经分析了很多容易出现竞争条件的例子。采用读-提交和快照隔离可以防止其中一部分，但并非对所有情况都有效，例如写倾斜和幻读所导致的棘手问题。最后你会发现面临以下挑战：



- 隔离级别通常难以理解，而且不同的数据库的实现不尽一致（例如“可重复读取”的含义在各家数据库的差别很大）。
- 如果去检查应用层的代码，往往很难判断它在特定的隔离级别下是否安全，特别是对于大型应用系统，几乎无法预测所有可能并发情况。
- 同时，还缺乏好的工具来帮助检测竞争状况。理论上，静态分析可能有所帮助<sup>[26]</sup>，但更多的还只是学术研究缺乏实用性。测试并发性问题往往效率很低，一切取决于时机，它只有在特定的情景下才会出现，存在很大的不确定性。

而这些都不是新问题，自20世纪70年代引入弱隔离级别<sup>[2]</sup>以后，这种情况就一直存在。长久以来，研究人员给出的答案都很简单：采用可串行化隔离！

可串行化隔离通常被认为是最强的隔离级别。它保证即使事务可能会并行执行，但最终的结果与每次一个即串行执行结果相同。这意味着，如果事务在单独运行时表现正确，那么它们在并发运行时结果仍然正确，换句话说，数据库可以防止所有可能的竞争条件。

如果串行化隔离比其他各种弱隔离级别好得多，那么为什么没有广泛使用呢？要回答这个问题，我们需要看看可串行化究竟是什么，以及如何执行。目前大多数提供可串行化的数据库都使用了以下三种技术之一，我们将依次探讨：

- 严格按照串行顺序执行（参阅本章后面的“实际的串行执行”）。
- 两阶段锁定（参阅本章后面的“两阶段加锁”），几十年来这几乎是唯一可行的选择。
- 乐观并发控制技术，例如可串行化的快照隔离（参阅本章后面的“可串行化的快照隔离”）。

我们首先限定在单节点背景下讨论这些技术。在接下来的第9章，我们将其推广到分布式系统多节点场景。

## 实际串行执行

解决并发问题最直接的方法是避免并发：即在一个线程上按顺序方式每次只执行一个事务。这样我们完全回避了诸如检测、防止事务冲突等问题，其对应的隔离级别一定是严格串行化的。

看上去这是个很直白的想法，但数据库设计人员直到2007年前后才完全确信，采用单线程循环来执行事务是可行的<sup>[45]</sup>。如果多线程并发在过去的30年中一只被认为是提升



性能的关键，那么现在转向单线程执行，这意味着什么呢？

有以下两方面的进展促使我们重新做出思考：

- 内存越来越便宜，现在许多应用可以将整个活动数据集都加载到内存中（参阅第3章“将所有内容加载到内存”）。当事务所需的所有数据都在内存中时，事务的执行速度要比等待磁盘I/O快得多。
- 数据库设计人员意识到OLTP事务通常执行很快，只产生少量的读写操作（参阅第3章“事务处理或分析”）。相比之下，运行时间较长的分析查询则通常是只读的，可以在一致性快照（使用快照隔离）上运行，而不需要运行在串行主循环里。

VoltDB / H-Store、Redis和Datomic等采用串行方式执行事务<sup>[46-48]</sup>。单线程执行有时可能会比支持并发的系统效率更高，尤其是可以避免锁开销。但是，其吞吐量上限是单个CPU核的吞吐量。为了充分利用单线程，相比于传统形式，事务也需要做出相应调整。

## 采用存储过程封装事务

在数据库的早期应用阶段，采用事务机制是希望能囊括用户的所有操作序列。例如，预订机票涉及多个步骤（搜索路线，票价和可用座位，决定行程，在行程的某个航班上预订座位，输入乘客信息，最后是付款）。数据库设计者认为，如果整个过程是一个事务，那么就可以方便地原子化执行。

然而，人类做出决定并回应的速度通常比较慢。如果数据库事务总是需要等待来自用户的输入，同时还要支持潜在大量并发需求，那么系统大部分时间将处于空闲状态。这样数据库无法高效运行，所以几乎所有的OLTP应用程序都避免在事务执行中等待用户交互从而使事务非常简洁。对于Web，这意味着事务会在一个HTTP请求中提交，而不会跨越多个请求。新的事务往往需要开启新的HTTP请求。

即使把人为交互从关键路径移除掉，事务总体沿用的依然是交互式客户端/服务器风格，一次一个请求语句。应用程序来提交查询，读取结果，可能会根据前一个查询的结果来进行其他查询，依此类推。请求与结果在应用层代码（某台机器）和数据库服务器（另一台机器）之间来回交互。

对于这种交互式的事务处理，大量时间花费在应用程序与数据库之间的网络通信。如果不允许事务并发，而是一次仅处理一个，那么吞吐量非常低，数据库总是在等待应





用提交下一个请求。在这种类型的数据库中，为了获得足够的吞吐性能，需要能够同时处理多个事务。

出于这个原因，采用单线程串行执行的系统往往不支持交互式的多语句事务。应用程序必须提交整个事务代码作为存储过程打包发送到数据库。这之间的差异如图7-9所示。把事务所需的所有数据全部加载在内存中，使存储过程高效执行，而无需等待网络或磁盘I/O。

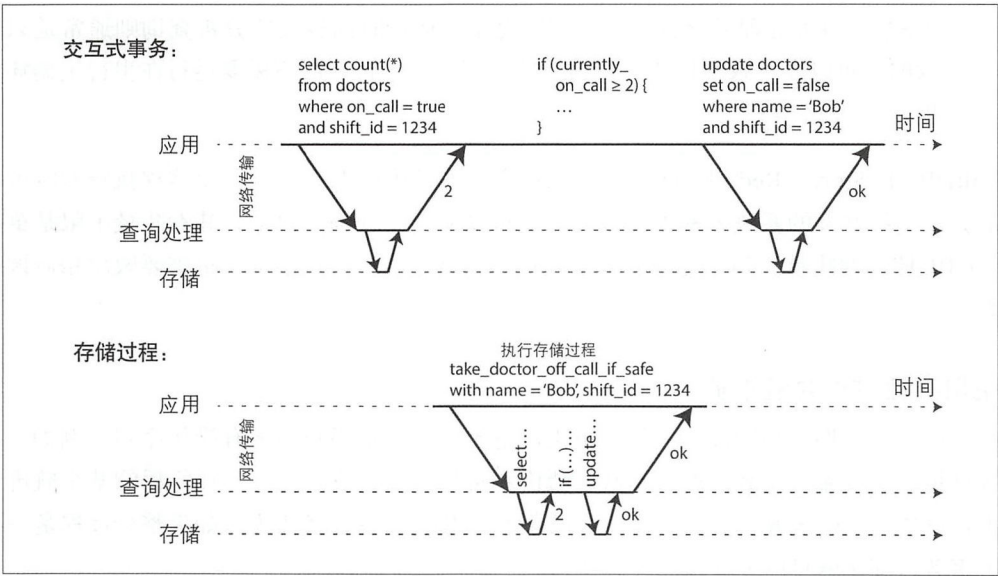


图7-9：交互式事务与存储过程的区别（使用图7-8的示例）

### 存储过程的优缺点

关系型数据库支持存储过程已经有很长一段时间了，自1999年以来已是SQL标准（SQL / PSM）。然而由于各种原因，存储过程的声誉有所下降：

- 每家数据库厂商都有自己的存储过程语言（Oracle的PL / SQL，SQL Server的T-SQL，PostgreSQL的PL / pgSQL等）。这些语言并没有跟上通用编程语言的发展，如果从今天的角度来看，这些语义都相当丑陋、过时，而且缺乏如今大多数编程语言所常用的函数库。
- 在数据库中运行代码难以管理：与应用服务器相比，调试更加困难，版本控制与部署复杂，测试不便，并且不容易和指标监控系统集成。
- 因为数据库实例往往被多个应用服务器所共享，所以数据库通常比应用服务器要



求更高的性能。数据库中一个设计不好的存储过程（例如，消耗大量内存或CPU时间）要比同样低效的应用服务器代码带来更大的麻烦。

不过这些问题也是可以克服的。最新的存储过程已经放弃了PL / SQL，而是使用现有的通用编程语言，例如VoltDB使用Java或Groovy，Datomic使用Java或Clojure，而Redis使用Lua。

存储过程与内存式数据存储使得单线程上执行所有事务变得可行。它们不需要等待I/O，避免加锁开销等复杂的并发控制机制，可以得到相当不错的吞吐量。

VoltDB还借助存储过程来执行复制：它并非将事务的执行结果从一个节点复制到另一个节点，而是在每个副本上都执行相同的存储过程。因此，VoltDB要求存储过程必须是确定性的（即不同的节点上运行时，结果必须完全相同）。如果事务需要获得当前的日期和时间等，必须通过专门的确定性API来实现。

## 分区

串行执行所有事务使得并发控制更加简单，但是数据库的吞吐量被限制在单机单个CPU核。虽然只读事务可以在单独的快照上执行，但是对于那些高写入需求的应用程序，单线程事务处理很容易成为严重的瓶颈。

为了扩展到多个CPU核和多节点，可以对数据进行分区（参见第6章），VoltDB支持这种配置模式。如果你能找到一个方法来对数据集进行分区，使得每个事务只在单个分区内读写数据，这样每个分区都可以有自己的事务处理线程且独立运行。此时为每个CPU核分配一个分区，则数据库的总体事务吞吐量可以到达与CPU核的数量成线性比例关系<sup>[47]</sup>。

但是，对于跨分区的事务，数据库必须在涉及的所有分区之间协调事务。存储过程需要跨越所有分区加锁执行，以确保整个系统的可串行化。

由于跨分区事务具有额外的协调开销，其性能比单分区内要慢得多。VoltDB报告的跨区事务吞吐量大约只有1000次/秒，这比单分区吞吐量低了好几个数量级，而且没法通过增加更多机器的方式来扩展性能<sup>[49]</sup>。

事务是否能只在单分区上执行很大程度上取决于应用层的数据结构。简单的键-值数据比较容易切分，而带有多个二级索引的数据则需要大量的跨区协调（参阅第6章“分区与二级索引”），因此不太合适。

## 串行执行小结

当满足以下约束条件时，串行执行事务可以实现串行化隔离：

- 事务必须简短而高效，否则一个缓慢的事务会影响到所有其他事务的执行性能。
- 仅限于活动数据集完全可以加载到内存的场景。有些很少访问的数据可能会被移到磁盘，但万一单线程事务需要访问它，就会严重拖累性能<sup>注10</sup>。
- 写入吞吐量必须足够低，才能在单个CPU核上处理；否则就需要采用分区，最好没有跨分区事务。
- 跨分区事务虽然也可以支持，但是占比必须很小。

## 两阶段加锁

近三十年来，可以说数据库只有一种被广泛使用的串行化算法，那就是两阶段加锁（two-phase locking, 2PL）<sup>注11</sup>。



### 2PL不是2PC

虽然两阶段加锁（2PL）听起来和两阶段提交（two-phase commit, 2PC）很相近，但它们是完全不同的东西。我们将在第9章讨论2PC。

之前我们看到，可以使用加锁的方法来防止脏写（参阅本章前面的“防止脏写”）：即如果两个事务同时尝试写入同一个对象时，以加锁的方式来确保第二个写入等待前面事务完成（包括中止或提交）。

两阶段加锁方法类似，但锁的强制性更高。多个事务可以同时读取同一对象，但只要出现任何写操作（包括修改或删除），则必须加锁以独占访问：

- 如果事务A已经读取了某个对象，此时事务B想要写入该对象，那么B必须等到A提交或中止之才能继续。以确保B不会在事务A执行的过程中去修改对象。
- 如果事务A已经修改了对象，此时事务B想要读取该对象，则B必须等到A提交或中止之后才能继续。对于2PL，不会出现读到旧值的情况（参见图7-1的示例）。

---

注10：如果事务需要访问那些不在内存的数据，最好的解决方案可能是中止事务，异步地将数据提取到内存中，同时继续处理其他事务，然后在数据加载完成后重启事务。这种方法被称为反高速缓存，在前面第2章有过描述。

注11：有时也被称为严格的两阶段加锁（strong strict two-phase locking, SS2PL），以区别于2PL的其他变体。



因此2PL不仅在并发写操作之间互斥，读取也会和修改产生互斥。快照级别隔离的口号“读写互不干扰”（参阅本章前面的“实现快照级别隔离”）非常准确地点明了它和两阶段加锁的关键区别。另一方面，因为2PL提供了串行化，所以它可以防止前面讨论的所有竞争条件，包括更新丢失和写倾斜。

## 实现两阶段加锁

目前，2PL已经用于MySQL（InnoDB）和SQL Server中的“可串行化隔离”，以及DB2中的“可重复读隔离”<sup>[23,36]</sup>。

此时数据库的每个对象都有一个读写锁来隔离读写操作。即锁可以处于共享模式或独占模式。基本用法如下：

- 如果事务要读取对象，必须先以共享模式获得锁。可以有多个事务同时获得一个对象的共享锁，但是如果某个事务已经获得了对应的独占锁，则所有其他事务必须等待。
- 如果事务要修改对象，必须以独占模式获取锁。不允许多个事务同时持有该锁（包括共享或独占模式），换言之，如果对象上已被加锁，则修改事务必须等待。
- 如果事务首先读取对象，然后尝试写入对象，则需要将共享锁升级为独占锁。升级锁的流程等价于直接获得独占锁。
- 事务获得锁之后，一直持有锁直到事务结束（包括提交或中止）。这也是名字“两阶段”的来由，在第一阶段即事务执行之前要获取锁，第二阶段（即事务结束时）则释放锁。

由于使用了这么多的锁机制，所以很容易出现死锁现象，例如事务A可能在等待事务B释放它的锁，而事务B在等待事务A释放所持有的锁。数据库系统会自动检测事务之间的死锁情况，并强行中止其中的一个以打破僵局，这样另一个可以继续向前执行。而被中止的事务需要由应用层来重试。

## 两阶段加锁的性能

两阶段加锁的主要缺点，或者说自1970年以来并不被所有人接纳的主要原因在于性能：其事务吞吐量和查询响应时间相比于其他弱隔离级别下降非常多。

部分原因在于锁的获取和释放本身的开销，但更重要的是其降低了事务的并发性。按2PL的设计，两个并发事务如果试图做任何可能导致竞争条件的事情，其中一个必须等待对方完成。

传统的关系数据库并不限制事务的执行时间，且当初是为和人类输入等交互式应用而设计的。结合2PL，最终结果是，当一个事务还需要等待另一个事务时，那么最终的等待时间几乎是没有限制的。即使可以保证自己的事务足够简短、高效，但一旦出现多个事务同时访问同一对象，会形成一个等待队列，事务就必须等待队列前面所有其他事务完成之后才能继续。

因此，在2PL模式下数据库的访问延迟具有非常大的不确定性，如果工作负载存在严重竞争，以百分比方式观察延迟指标会发现非常缓慢（参阅第1章“描述性能”）。试想这种情况，某个事务本身很慢，或者是由于需要访问大量数据而获得了许多锁，则它还会导致系统的其他部分都停顿下来。如果应用需要稳定如一的性能，这种不确定性就是致命的。

同样是基于加锁方式的读-提交隔离也可能发生死锁，然而在2PL下，取决于事务的访问模式，死锁可能变得更为频繁。因而导致另一个性能问题，即如果事务由于死锁而被强行中止，应用层就必须从头重试，假如死锁过于频繁，则最后的性能和效率必然大打折扣。

## 谓词锁

对于加锁，我们还忽略了一个微妙但重要的细节。如本章前面“写倾斜与幻读”中的幻读问题，即一个事务改变另一个事务的查询结果。可串行化隔离也必须防止幻读问题。

以会议室预订为例，如果事务在查询某个时间段内一个房间的预订情况（参见示例7-2），则另一个事务不能同时去插入或更新同一时间段内该房间的预订情况，但它修改其他房间的预订情况，或者在不影响当前查询的情况下，修改该房间的其他时间段预订。

如何实现呢？技术上讲，我们需要引入一种谓词锁（或者属性谓词锁）<sup>[3]</sup>。它的作用类似于之前描述的共享/独占锁，而区别在于，它并不属于某个特定的对象（如表的某一行），而是作用于满足某些搜索条件的所有查询对象，例如：

```
SELECT * FROM bookings
WHERE room_id = 123 AND
      end_time > '2018-01-01 12:00' AND
      start_time < '2018-01-01 13:00';
```

谓词锁会限制如下访问：

- 如果事务A想要读取某些满足匹配条件的对象，例如采用SELECT查询，它必须以



共享模式获得查询条件的谓词锁。如果另一个事务B正持有任何一个匹配对象的互斥锁，那么A必须等到B释放锁之后才能继续执行查询。

- 如果事务A想要插入、更新或删除任何对象，则必须首先检查所有旧值和新值是否与现有的任何谓词锁匹配（即冲突）。如果事务B持有这样的谓词锁，那么A必须等到B完成提交（或中止）后才能继续。

这里的关键点在于，谓词锁甚至可以保护数据库中那些尚不存在但可能马上会被插入的对象（幻读）。将两阶段加锁与谓词锁结合使用，数据库可以防止所有形式的写倾斜以及其他竞争条件，隔离变得真正可串行化。

## 索引区间锁

不幸的是，谓词锁性能不佳：如果活动事务中存在许多锁，那么检查匹配这些锁就变得非常耗时。因此，大多数使用2PL的数据库实际上实现的是索引区间锁（或者*next-key locking*），本质上它是对谓词锁的简化或者近似<sup>[41,50]</sup>。

简化谓词锁的方式是将其保护的對象扩大化，首先这肯定是安全的。例如，如果一个谓词锁保护的是查询条件是：房间123，时间段是中午至下午1点，则一种方式是通过扩大时间段来简化，即保护123房间的所有时间段；或者另一种方式是扩大房间，即保护中午至下午1点之间的所有房间（而不仅是123号房间）。这样，任何与原始谓词锁冲突的操作肯定也和近似之后的区间锁相冲突。

对于房间预订数据库，通常会在`room_id`列上创建索引，和/或在`start_time`和`end_time`上有索引（否则前面的查询在大型数据库上会很慢）：

- 假设索引位于`room_id`上，数据库使用此索引查找123号房间的当前预订情况。现在，数据库可以简单地将共享锁附加到此索引条目，表明事务已搜索了123号房间的所有时间段预订。
- 或者，如果数据库使用基于时间的索引来查找预订，则可以将共享锁附加到该索引中的一系列值，表示事务已经搜索了该时间段内的所有值（例如直到2020年1月1日）。

无论哪种方式，查询条件的近似值都附加到某个索引上。接下来，如果另一个事务想要插入、更新或删除同一个房间和/或重叠时间段的预订，则肯定需要更新这些索引，一定就会与共享锁冲突，因此会自动处于等待状态直到共享锁释放。

这样就有效防止了写倾斜和幻读问题。的确，索引区间锁不像谓词锁那么精确（会锁





定更大范围的对象，而超出了串行化所要求的部分），但由于开销低得多，可以认为是一种很好的折衷方案。

如果没有合适的索引可以施加区间锁，则数据库可以回退到对整个表施加共享锁。这种方式的性能肯定不好，它甚至会阻止所有其他事务的写操作，但的确可以保证安全性。

## 可串行化的快照隔离

本章已经给大家展示了数据库并发方面很多让人纠结、黯淡的一面。两阶段加锁虽然可以保证串行化，但性能差强人意且无法扩展（由于串行执行）；弱级别隔离虽然性能不错，但容易引发各种边界条件（如更新丢失，写倾斜，幻读等）。那么，串行化隔离与性能是不是从根本上就是互相冲突而无法兼得吗？

或许并非如此。最近一种称为可串行化的快照隔离（Serializable Snapshot Isolation, SSI）算法看起来让人眼前一亮。它提供了完整的可串行性保证，而性能相比于快照隔离损失很小。SSI算法面世至今不过数年，它于2008年被首次提出<sup>[40]</sup>，后来成为Michael Cahill的博士论文研究主题<sup>[51]</sup>。

目前，SSI可用于单节点数据库（PostgreSQL 9.1之后的可串行化隔离）或者分布式数据库（如FoundationDB采用了类似的算法）。相比于其他并发控制机制，SSI尚需在实践中证明其性能。即使如此，它很有可能成为未来数据库的标配。

## 悲观与乐观的并发控制

两阶段加锁是一种典型的悲观并发控制机制。它基于这样的设计原则：如果某些操作可能出错（例如与其他并发事务发生了锁冲突），那么直接放弃，采用等待方式直到绝对安全。这和多线程编程中互斥锁是一致的。

某种意义上讲，串行执行是种极端悲观的选择：事务执行期间，等价于事务对整个数据库（或数据库的一个分区）持有互斥锁。而我们只能假定事务执行得足够快、持锁时间足够短，来稍稍弥补这种悲观色彩。

相比之下，可串行化的快照隔离则是一种乐观并发控制。在这种情况下，如果可能发生潜在冲突，事务会继续执行而不是中止，寄希望一切相安无事；而当事务提交时（只有可串行化的事务被允许提交），数据库会检查是否确实发生了冲突（即违反了隔离性原则），如果是的话，中止事务并接下来重试。

乐观并发控制其实是一个古老的想法<sup>[52]</sup>，关于其优点和缺点已经争论了很长时



间<sup>[53]</sup>。如果冲突很多，则性能不佳（许多事务试图访问相同的对象），大量的事务必须中止。如果系统已接近其最大吞吐量，反复重试事务会使系统性能变得更差。

但是，如果系统还有足够的性能提升空间，且如果事务之间的竞争不大，乐观并发控制会比悲观方式高效很多。通过可交换的原子操作还可以减少一些竞争情况。例如，如果多个事务同时试图增加某个计数器，那么不管以什么样的顺序去增加（只要同一事务不去读计数器），最后的结果总是等价的，并发提交多个增量操作是可行的。

顾名思义，SSI基于快照隔离，也就是说，事务中的所有读取操作都是基于数据库的一致性快照（请参阅本章前面的“快照隔离”和“可重复读”）。这是与早期的乐观并发控制主要区别。在快照隔离的基础上，SSI新增加了相关算法来检测写入之间的串行化冲突从而决定中止哪些事务。

## 基于过期的条件做决定

我们在讨论写倾斜（参阅本章前面的“写倾斜与幻读”）时，介绍了这样一种使用场景：事务首先查询某些数据，根据查询的结果来决定采取后续操作，例如修改数据。而在快照隔离情况下，数据可能在查询期间就已经被其他事务修改，导致原事务在提交时决策的依据信息已出现变化。

换句话说，事务是基于某些前提条件而决定采取行动，在事务开始时条件成立，例如“目前有两名医生值班”，而当事务要提交时，数据可能已经发生改变，条件已不再成立。

当应用程序执行查询时（例如“当前有多少医生在值班？”），数据库本身无法预知应用层逻辑如何使用这些查询结果。安全起见，数据库假定对查询结果（决策的前提条件）的任何变化都应使写事务失效。换言之，查询与写事务之间可能存在因果依赖关系。为了提供可串行化的隔离，数据库必须检测事务是否会修改其他事务的查询结果，并在此情况下中止写事务。

数据库如何知道查询结果是否发生了改变呢？可以分以下两种情况：

- 读取是否作用于一个（即将）过期的MVCC对象（读取之前已经有未提交的写入）。
- 检查写入是否影响即将完成的读取（读取之后，又有新的写入）。

## 检测是否读取了过期的MVCC对象

回想一下，快照隔离通常采用多版本并发控制技术（MVCC，见图7-10）来实现。当



事务从MVCC数据库一致性快照读取时，它会忽略那些在创建快照时尚未提交的事务写入。例如图7-10中，事务42（修改Alice的值班状态）未被提交，因此事务43中Alice查询到的on\_call是true；当事务43提交时，事务42已经完成了提交。换言之，从快照读取时被忽略的写入已经生效，并且直接导致事务43做决定的前提已不再成立。

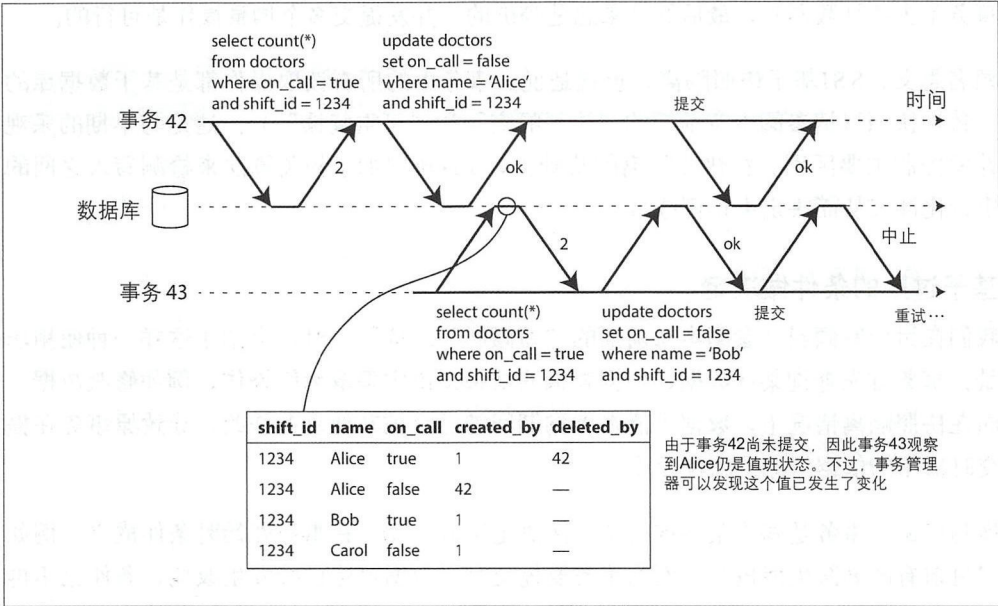


图7-10：检测事务是否从MVCC快照中读取了旧值

为防止这种异常，数据库需要跟踪那些由于MVCC可见性规则而被忽略的写操作。当事务提交时，数据库会检查是否存在一些当初被忽略的写操作现在已经完成了提交，如果是则必须中止当前事务。

为什么要等到提交：当检测到读旧值，为何不立即中止事务43呢？可以考虑这些情况，首先，如果事务43是个只读事务，没有任何写倾斜风险，就不需要中止；而事务43读取数据库时，数据库还不知道事务是否稍后有任何写操作。此外，事务43提交时，有可能事务42发生了中止或者还处于未提交状态，因此读取的并非是过期值。通过减少不必要的中止，SSI可以高效支持那些需要在一致性快照中运行很长时间的读事务。

## 检测写是否影响了之前的读

第二种要考虑的情况是，在读取数据之后，另一个事务修改了数据。如图7-11所示。



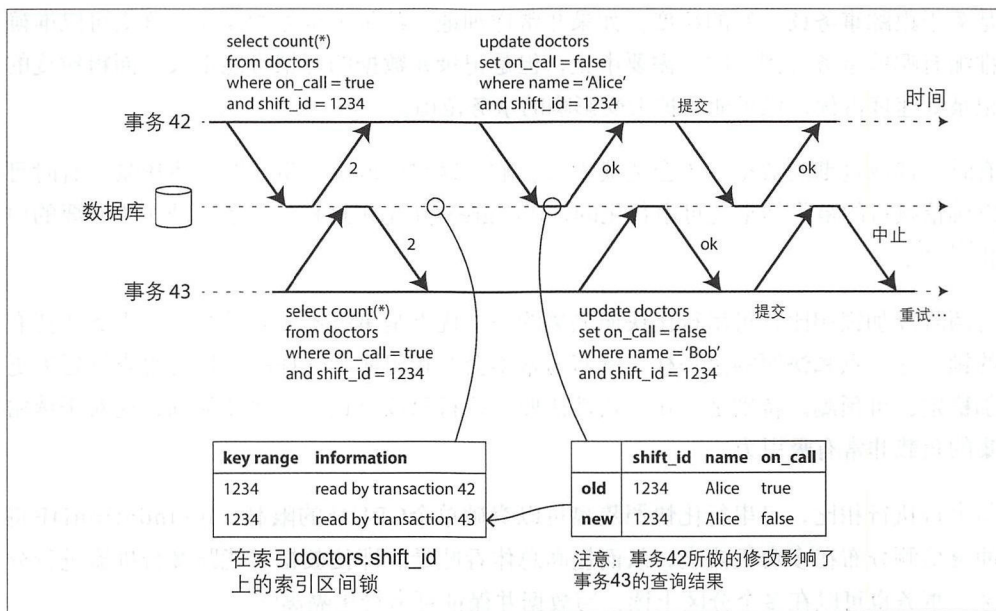


图7-11：在可串行化的快照隔离中，检测事务是否修改了另一个事务的查询结果

在“两阶段加锁”中，我们讨论了索引区间锁（参阅本章前面的“索引区间锁”），它可以锁定与某个查询条件匹配的所有行，例如 `WHERE shift_id = 1234`。这里使用了类似的技术，只有一点差异：SSI锁不会阻塞其他事务。

在图7-11中，事务42和事务43都在查询轮班1234期间的值班医生。如果在 `shift_id` 上建有索引，数据库可以通过索引条目1234来记录事务42和事务43都查询了相同的结果。如果没有索引，可以在表级别跟踪此信息。该额外记录只需保留很小一段时间，当并发的所有事务都处理完成（提交或中止）之后，就可以丢弃。

当另一个事务尝试修改时，它首先检查索引，从而确定是否最近存在一些读目标数据的其他事务。这个过程类似于在受影响的字段范围上获取写锁，但它并不会阻塞读取，而是直到读事务提交时才进一步通知他们：所读到的数据现在已经发生了变化。

图7-11中，事务43和事务42会互相通知对方先前的读已经过期。虽然事务43的修改的确影响了事务42，但事务43当时并未提交（修改未生效），而事务42首先尝试提交，所以可以成功；随后当事务43试图提交时，来自42的冲突写已经提交生效，事务43不得不中止。

## 可串行化快照隔离的性能

有许多工程方面的细节会直接影响算法在实践中的效果。例如，一个需要权衡考虑的

是关于跟踪事务读、写的粒度。如果非常详细地跟踪每个事务的操作，确实可以准确推测有哪些事务受到影响、需要中止，但是记录元数据的开销可能很大；而粗粒度的记录则速度占优，但可能会扩大受影响的事务范围。

有时，读取过期的数据并不会造成太大影响，这完全取决于所处的具体场景。有时可以确信执行的最终结果是可串行化的，PostgreSQL采用这样的信条来减少不必要的中止<sup>[11,41]</sup>。

与两阶段加锁相比，可串行化快照隔离的一大优点是事务不需要等待其他事务所持有的锁。这一点和快照隔离一样，读写通常不会互相阻塞。这样的设计使得查询延迟更加稳定、可预测。特别是，在一致性快照上执行只读查询不需要任何锁，这对于读密集的负载非常有吸引力。

与串行执行相比，可串行化快照隔离可以突破单个CPU核的限制。FoundationDB将冲突检测分布在多台机器上，从而提高总体吞吐量。即使数据可能跨多台机器进行分区，事务也可以在多个分区上读、写数据并保证可串行化隔离<sup>[54]</sup>。

需要指出，事务中止的比例会显著影响SSI的性能表现。例如，一个运行很长时间的事务，读取和写入了大量数据，因而产生冲突并中止的概率就会增大，所以SSI要求读-写型事务要简短（而长时间执行的只读事务则没有此限制）。但总体讲，相比于两阶段加锁与串行执行，SSI更能容忍那些执行缓慢的事务。

## 小结

事务作为一个抽象层，使得应用程序可以忽略数据库内部一些复杂的并发问题，以及某些硬件、软件故障，从而简化应用层的处理逻辑，大量的错误可以转化为简单的事务中止和应用层重试。

本章，我们例举了很多事务能够预防的问题。尽管并非所有的应用程序都会面临这些问题，例如那些简单的访问模式，只读或者只写，可能根本无需事务的帮助。但对于复杂的访问模式，事务可以大大简化需要考虑的潜在错误情况。

如果没有事务，各种错误情况（如进程崩溃、网络中断、停电、磁盘已满、并发竞争等）会导致数据可能出现各种不一致。例如，反规格化数据模式和相关操作会导致与源数据不同步。假如没有事务，处理那些复杂交互访问最后所导致的数据库混乱就会异常痛苦。

本章，我们深入探讨了并发控制这一主题。介绍了多个广泛使用的隔离级别，特别是



读-提交，快照隔离（或可重复读取）与可串行化。通过分析如何处理边界条件来阐述这些隔离级别的要点：

### 脏读

客户端读到了其他客户端尚未提交的写入。读-提交以及更强的隔离级别可以防止脏读。

### 脏写

客户端覆盖了另一个客户端尚未提交的写入。几乎所有的数据库实现都可以防止脏写。

### 读倾斜（不可重复读）

客户在不同的时间点看到了不同值。快照隔离是最用的防范手段，即事务总是在某个时间点的一致性快照中读取数据。通常采用多版本并发控制（MVCC）来实现快照隔离。

### 更新丢失

两个客户端同时执行读-修改-写入操作序列，出现了其中一个覆盖了另一个的写入，但又没有包含对方最新值的情况，最终导致了部分修改数据发生了丢失。快照隔离的一些实现可以自动防止这种异常，而另一些则需要手动锁定查询结果（SELECT FOR UPDATE）。

### 写倾斜

事务首先查询数据，根据返回的结果而作出某些决定，然后修改数据库。当事务提交时，支持决定的前提条件已不再成立。只有可串行化的隔离才能防止这种异常。

### 幻读

事务读取了某些符合查询条件的对象，同时另一个客户端执行写入，改变了先前的查询结果。快照隔离可以防止简单的幻读，但写倾斜情况则需要特殊处理，例如采用区间范围锁。

弱隔离级别可以防止上面的某些异常，但还需要应用开发人员手动处理其他复杂情况（例如，显式加锁）。只有可串行化的隔离可以防止所有这些问题。我们主要讨论了实现可串行化隔离的三种不同方法：

### 严格串行执行事务

如果每个事务的执行速度非常快，且单个CPU核可以满足事务的吞吐量要求，严格串行执行是一个非常简单有效的方案。





## 两阶段加锁

几十年来，这一直是实现可串行化的标准方式，但还是有很多系统出于性能原因而放弃使用它。

## 可串行化的快照隔离 (SSI)

一种最新的算法，可以避免前面方法的大部分缺点。它秉持乐观预期的原则，允许多个事务并发执行而不互相阻塞；仅当事务尝试提交时，才检查可能的冲突，如果发现违背了串行化，则某些事务会被中止。

本章中的示例都采用关系数据模型。但是，正如本章前面的“多对象事务的需求”所描述的，无论对哪种数据模型，事务都是非常有用的数据库功能。

本章所介绍的算法、方案主要针对单节点。对于分布式数据库，还会面临更多、更复杂的挑战，我们将在接下来的两章中继续展开讨论。

## 参考文献

- [1] Donald D. Chamberlin, Morton M. Astrahan, Michael W. Blasgen, et al.: “A History and Evaluation of System R,” *Communications of the ACM*, volume 24, number 10, pages 632-646, October 1981. doi:10.1145/358769.358784.
- [2] Jim N. Gray, Raymond A. Lorie, Gianfranco R. Putzolu, and Irving L. Traiger: “Granularity of Locks and Degrees of Consistency in a Shared Data Base,” in *Modelling in Data Base Management Systems: Proceedings of the IFIP Working Conference on Modelling in Data Base Management Systems*, edited by G. M. Nijssen, pages 364 – 394, Elsevier/North Holland Publishing, 1976. Also in *Readings in Database Systems*, 4th edition, edited by Joseph M. Hellerstein and Michael Stonebraker, MIT Press, 2005. ISBN: 978-0-262-69314-1
- [3] Kapali P. Eswaran, Jim N. Gray, Raymond A. Lorie, and Irving L. Traiger: “The Notions of Consistency and Predicate Locks in a Database System,” *Communications of the ACM*, volume 19, number 11, pages 624-633, November 1976.
- [4] “ACID Transactions Are Incredibly Helpful,” FoundationDB, LLC, 2013.
- [5] John D. Cook: “ACID Versus BASE for Database Transactions,” *johnndcook.com*, July 6, 2009.
- [6] Gavin Clarke: “NoSQL’s CAP Theorem Busters: We Don’t Drop ACID,”

*theregister.co.uk*, November 22, 2012.

[7] Theo Härder and Andreas Reuter: “Principles of Transaction-Oriented Database Recovery,” *ACM Computing Surveys*, volume 15, number 4, pages 287-317, December 1983. doi:10.1145/289.291.

[8] Peter Bailis, Alan Fekete, Ali Ghodsi, et al.: “HAT, not CAP: Towards Highly Available Transactions,” at *14th USENIX Workshop on Hot Topics in Operating Systems (HotOS)*, May 2013.

[9] Armando Fox, Steven D. Gribble, Yatin Chawathe, et al.: “Cluster-Based Scalable Network Services,” at *16th ACM Symposium on Operating Systems Principles (SOSP)*, October 1997.

[10] Philip A. Bernstein, Vassos Hadzilacos, and Nathan Goodman: *Concurrency Control and Recovery in Database Systems*. Addison-Wesley, 1987. ISBN: 978-0-201-10715-9, available online at *research.microsoft.com*.

[11] Alan Fekete, Dimitrios Liarokapis, Elizabeth O’Neil, et al.: “Making Snapshot Isolation Serializable,” *ACM Transactions on Database Systems*, volume 30, number 2, pages 492-528, June 2005. doi:10.1145/1071610.1071615.

[12] Mai Zheng, Joseph Tucek, Feng Qin, and Mark Lillibridge: “Understanding the Robustness of SSDs Under Power Fault,” at *11th USENIX Conference on File and Storage Technologies (FAST)*, February 2013.

[13] Laurie Denness: “SSDs: A Gift and a Curse,” *laur.ie*, June 2, 2015.

[14] Adam Surak: “When Solid State Drives Are Not That Solid,” *blog.algolia.com*, June 15, 2015.

[15] Thanumalayan Sankaranarayanan Pillai, Vijay Chidambaram, Ramnathan Alagappan, et al.: “All File Systems Are Not Created Equal: On the Complexity of Crafting Crash-Consistent Applications,” at *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, October 2014.

[16] Chris Siebenmann: “Unix’s File Durability Problem,” *utcc.utoronto.ca*, April 14, 2016.

[17] Lakshmi N. Bairavasundaram, Garth R. Goodson, Bianca Schroeder, et al.: “An

Analysis of Data Corruption in the Storage Stack,” at *6th USENIX Conference on File and Storage Technologies* (FAST), February 2008.

[18] Bianca Schroeder, Raghav Lagisetty, and Arif Merchant: “Flash Reliability in Production: The Expected and the Unexpected,” at *14th USENIX Conference on File and Storage Technologies* (FAST), February 2016.

[19] Don Allison: “SSD Storage - Ignorance of Technology Is No Excuse,” *blog.korelogic.com*, March 24, 2015.

[20] Dave Scherer: “Those Are Not Transactions (Cassandra 2.0),” *blog.foundationdb.com*, September 6, 2013.

[21] Kyle Kingsbury: “Call Me Maybe: Cassandra,” *aphyr.com*, September 24, 2013.

[22] “ACID Support in Aerospike,” Aerospike, Inc., June 2014.

[23] Martin Kleppmann: “Hermitage: Testing the ‘I’ in ACID,” *martin.kleppmann.com*, November 25, 2014.

[24] Tristan D’Agosta: “BTC Stolen from Poloniex,” *bitcointalk.org*, March 4, 2014.

[25] bitcointhief2: “How I Stole Roughly 100 BTC from an Exchange and How I Could Have Stolen More!,” *reddit.com*, February 2, 2014.

[26] Sudhir Jorwekar, Alan Fekete, Krithi Ramamritham, and S. Sudarshan: “Automating the Detection of Snapshot Isolation Anomalies,” at *33rd International Conference on Very Large Data Bases* (VLDB), September 2007.

[27] Michael Melanson: “Transactions: The Limits of Isolation,” *michaelmelanson.net*, March 20, 2014.

[28] Hal Berenson, Philip A. Bernstein, Jim N. Gray, et al.: “A Critique of ANSI SQL Isolation Levels,” at *ACM International Conference on Management of Data* (SIGMOD), May 1995.

[29] Atul Adya: “Weak Consistency: A Generalized Theory and Optimistic Implementations for Distributed Transactions,” PhD Thesis, Massachusetts Institute of Technology, March 1999.

[30] Peter Bailis, Aaron Davidson, Alan Fekete, et al.: “Highly Available Transactions:



Virtues and Limitations (Extended Version),” at *40th International Conference on Very Large Data Bases (VLDB)*, September 2014.

[31] Bruce Momjian: “MVCC Unmasked,” *momjian.us*, July 2014.

[32] Annamalai Gurusami: “Repeatable Read Isolation Level in InnoDB - How Consistent Read View Works,” *blogs.oracle.com*, January 15, 2013.

[33] Nikita Prokopov: “Unofficial Guide to Datomic Internals,” *tonsky.me*, May 6, 2014.

[34] Baron Schwartz: “Immutability, MVCC, and Garbage Collection,” *xaprb.com*, December 28, 2013.

[35] J. Chris Anderson, Jan Lehnardt, and Noah Slater: *CouchDB: The Definitive Guide*. O'Reilly Media, 2010. ISBN: 978-0-596-15589-6.

[36] Rikdeb Mukherjee: “Isolation in DB2 (Repeatable Read, Read Stability, Cursor Stability, Uncommitted Read) with Examples,” *mframes.blogspot.co.uk*, July 4, 2013.

[37] Steve Hilker: “Cursor Stability (CS) - IBM DB2 Community,” *toadworld.com*, March 14, 2013.

[38] Nate Wiger: “An Atomic Rant,” *nateware.com*, February 18, 2010.

[39] Joel Jacobson: “Riak 2.0: Data Types,” *blog.joeljacobsen.com*, March 23, 2014.

[40] Michael J. Cahill, Uwe Röhm, and Alan Fekete: “Serializable Isolation for Snapshot Databases,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2008. doi:10.1145/1376616.1376690.

[41] Dan R. K. Ports and Kevin Grittner: “Serializable Snapshot Isolation in PostgreSQL,” at *38th International Conference on Very Large Databases (VLDB)*, August 2012.

[42] Tony Andrews: “Enforcing Complex Constraints in Oracle,” *tonyandrews.blogspot.co.uk*, October 15, 2004.

[43] Douglas B. Terry, Marvin M. Theimer, Karin Petersen, et al.: “Managing Update Conflicts in Bayou, a Weakly Connected Replicated Storage System,” at *15th ACM Symposium on Operating Systems Principles (SOSP)*, December 1995. doi:

10.1145/224056.224070.

[44] Gary Fredericks: “Postgres Serializability Bug,” *github.com*, September 2015.

[45] Michael Stonebraker, Samuel Madden, Daniel J. Abadi, et al.: “The End of an Architectural Era (It’s Time for a Complete Rewrite),” at *33rd International Conference on Very Large Data Bases (VLDB)*, September 2007.

[46] John Hugg: “H-Store/VoltDB Architecture vs. CEP Systems and Newer Streaming Architectures,” at *Data @Scale Boston*, November 2014.

[47] Robert Kallman, Hideaki Kimura, Jonathan Natkins, et al.: “H-Store: A High-Performance, Distributed Main Memory Transaction Processing System,” *Proceedings of the VLDB Endowment*, volume 1, number 2, pages 1496-1499, August 2008.

[48] Rich Hickey: “The Architecture of Datomic,” *infoq.com*, November 2, 2012.

[49] John Hugg: “Debunking Myths About the VoltDB In-Memory Database,” *voltldb.com*, May 12, 2014.

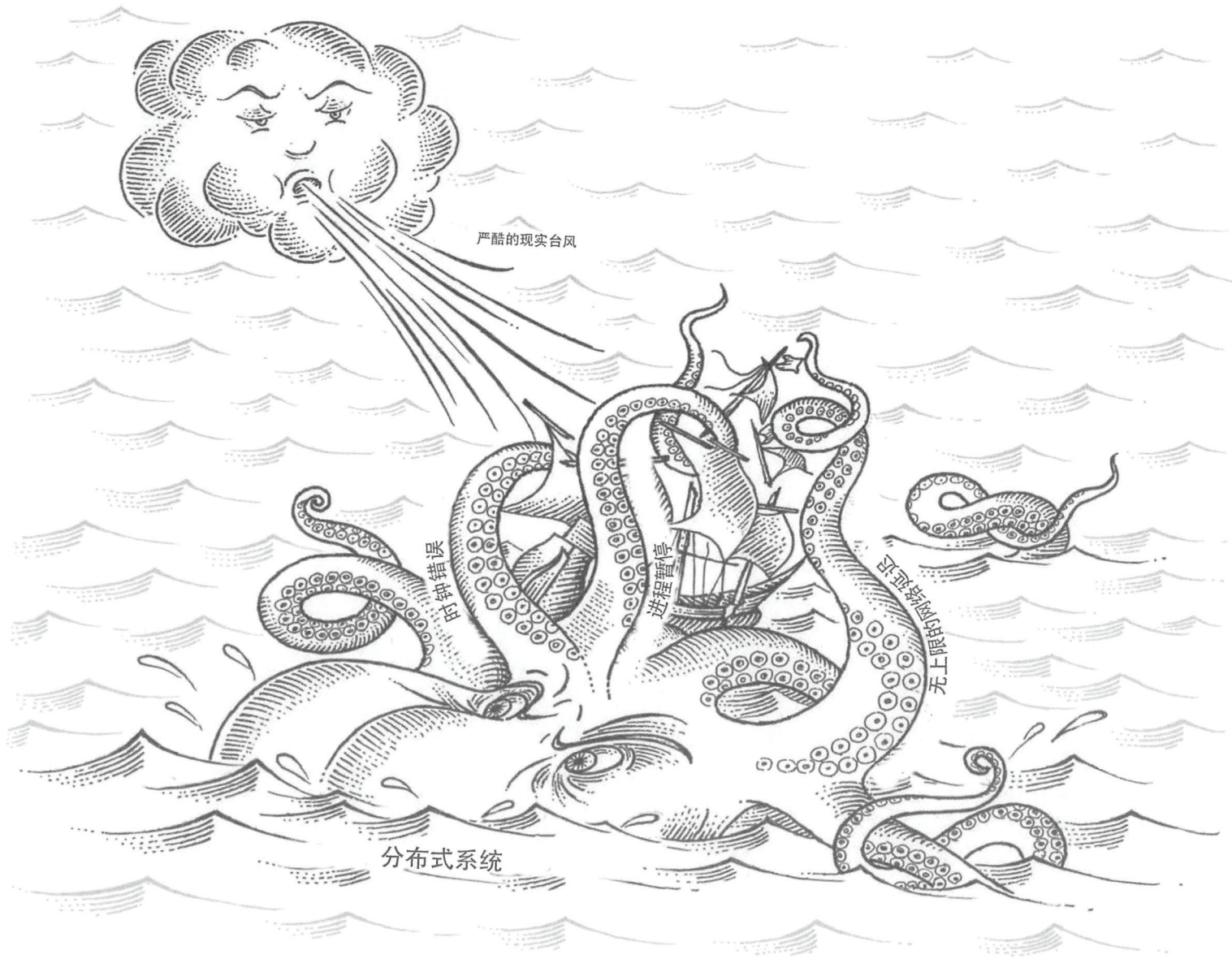
[50] Joseph M. Hellerstein, Michael Stonebraker, and James Hamilton: “Architecture of a Database System,” *Foundations and Trends in Databases*, volume 1, number 2, pages 141-259, November 2007. doi:10.1561/19000000002.

[51] Michael J. Cahill: “Serializable Isolation for Snapshot Databases,” PhD Thesis, University of Sydney, July 2009.

[52] D. Z. Badal: “Correctness of Concurrency Control and Implications in Distributed Databases,” at *3rd International IEEE Computer Software and Applications Conference (COMPSAC)*, November 1979.

[53] Rakesh Agrawal, Michael J. Carey, and Miron Livny: “Concurrency Control Performance Modeling: Alternatives and Implications,” *ACM Transactions on Database Systems (TODS)*, volume 12, number 4, pages 609-654, December 1987. doi: 10.1145/32204.32220.

[54] Dave Rosenthal: “Databases at 14.4MHz,” *blog.foundationdb.com*, December 10, 2014.





# 分布式系统的挑战

嘿，我刚刚认识你，  
网络滞后问题，  
这里是我的数据，  
希望你能保存它们。

——Kyle Kingsbury, Carly Rae Jepsen 与网络分区风险专栏系列 (2013)

前面几章反复出现的主题是系统出错之后该如何处理。例如，我们已经讨论了多副本切换（第5章的“处理节点失效”），复制滞后（第5章的“复制滞后问题”）以及事务的并发控制（第7章的“弱隔离级别”）等。只有更充分地了解实际系统可能出现的各种临界条件，我们才能更好地处理它们。

虽然长篇累牍谈了很多，但是这几章所剖析的错误情况还是过于乐观，现实则更加复杂。现在我们将这种复杂性推向极致，并做一个非常悲观的假定：所有可能出错的事情一定会出错<sup>注1</sup>。实际上，有经验的系统管理员会告诉你这其实是一个非常合理的假定，他们可能都遭遇过一些可怕甚至惨痛的教训。

关键在于，分布式系统与在单节点上的软件有着非常显著的区别，你会碰到五花八门、千奇百怪的问题所导致的各种故障<sup>[1,2]</sup>。本章我们将了解这些实践中所出现的问题，探明哪些称得上可靠而哪些是不可靠。

作为开发者，我们的核心任务是构建可靠的系统，即使系统面临各种出错可能，也需要完成预定工作（确保满足用户期望）。在下一章，我们将介绍分布式环境下可以提

注1： 此处暂时不考虑拜占庭式的错误（详见本章后面的“拜占庭故障”）。

供这类保证的常用算法和具体例子。但首先，本章的主要任务还是充分认识眼前的挑战。

本章对分布式系统可能出现的故障做了一个全面、近乎悲观的总结。故障可能来自网络问题（详见本章后面的“不可靠的网络”），以及时钟与时序问题（详见本章后面的“不可靠的时钟”）等，并讨论这些问题的可控程度。坦白讲，这些问题一直在困扰着人们，因此我们将探讨如何认清分布式系统的状态本质，并据此来评估所发生的各种故障（详见本章后面的“知识，真相与谎言”）。

## 故障与部分失效

在单节点上开发程序时，通常它应该以一种确定性的方式运行：要么工作，要么出错。有bug的软件可能最终会在某天暴露出来（而问题通常可以由重启机器来解决），但问题的主要原因可能还是软件本身的bug。

单台节点上的软件通常不应该出现模棱两可的现象：当硬件正常工作时，相同的操作通常总会产生相同的结果（即确定性）；而如果硬件存在问题（例如，内存损坏或接口松动），结果往往是系统性故障，如内核崩溃，蓝屏死机，启动失败等。因此在单节点上一个质量合格的软件状态要么是功能正常，要么是完全失效，而不会介于两者之间。

这背后涉及计算机设计一个非常审慎的选择：如果发生了某种内部错误，我们宁愿使计算机全部崩溃，而不是返回一个错误的结果，错误的结果往往更难处理。因此，计算机隐藏了一切模糊的物理世界，呈现以一个理想化的系统模型，像以数学完美的方式运行。CPU指令通常以确定性方式操作，如果写入一些数据到内存或磁盘，那么这些数据将保持不变且不会被随机破坏。这种确定-正确性的设计目标可以一直追溯到第一台数字计算机<sup>[3]</sup>。

然而当涉及多台节点时，情况发生了根本性变化。对于这种分布式系统，理想化的标准正确模型不再适用，我们必须面对一个可能非常混乱的现实。正如博客<sup>[4]</sup>所描述的那样，在这样一个现实世界中，各种各样的事情都可能会出错：

在本人的从业经历中，处理过数据中心网络分区故障，PDU（配电单元）故障，交换机故障，机架的意外断电，数据中心主干网故障，数据中心电力故障，以及某个司机将其福特皮卡撞向暖通空调等。另外，我得声明我其实并不是一个运维人员。

—— Coda Hale



在分布式系统中，可能会出现系统的一部分工作正常，但其他某些部分出现难以预测的故障，我们称之为“部分失效”。问题的难点就在于这种部分失效是不确定的：如果涉及多个节点和网络，几乎肯定会碰到有时网络正常，有时则莫名的失败。正如接下来马上要看到的，通过网络发送消息的延迟非常不确定，有时甚至根本不知道执行是否成功。

正是由于这种不确定性和部分失效大大提高了分布式系统的复杂性<sup>[5]</sup>。

## 云计算和超算

关于如何构建大规模计算系统有以下几种不同的思路：

- 规模的一个极端是高性能计算（high-performance computing, HPC）。包含成千上万个CPU的超级计算机构成一个庞大的集群，通常用于计算密集型的科学计算任务，如天气预报或分子动力学（模拟原子和分子的运动）。
- 另一个极端是云计算。虽然云计算的定义并非那么明确，但通常它具有以下特征：多租户数据中心，通用计算机，用IP以太网链接，弹性/按需资源分配，并按需计费。
- 传统企业数据中心则位于以上两个极端之间。

不同的集群构建方式所对应的错误处理方法也不尽相同。对于高性能计算，通常会定期对任务状态进行快照，然后保存在持久存储上，当某节点出现故障，解决方案是简单地停止整个集群的任务；等故障节点修复之后，从最近的快照检查点继续执行<sup>[7,8]</sup>。从这一点上看，高性能计算其实更像是一个单节点系统而不是分布式系统，它将局部失效升级为整体失效，例如系统的任何部分发生了故障，就干脆让系统停下来，这和单机上内核崩溃类似。

本书的重点是基于互联网的服务系统，这些系统与上述高性能计算有很多不同之处：

- 许多互联网服务都是在线的，需要随时（如 $7 \times 24\text{h}$ ）为用户提供低延迟服务。任何服务不可用情况，如停下集群来修复故障，都是不可取的。相比之下，像天气模拟这样的离线（批处理）作业则可以暂停然后重启，影响相对较小。
- 高性能计算通常采用专用硬件，每个节点的可靠性很高，节点间主要通过共享内存或者远程内存直接访问（RDMA）等技术进行通信。而云计算中的节点多是由通用机器构建，出于大规模部署时经济因素的考虑，单节点的成本低廉主要而依靠较高的集群聚合性能，但另一方面也具有较高的故障率。





- 大型数据中心网络通常基于IP和以太网，采用Clos拓扑结构提供等分带宽<sup>[9]</sup>。高性能计算则通常特定的网络拓扑结构，例如多维网格和toruses<sup>[10]</sup>，它们可以为HPC特定工作负载提供了更好的性能。
- 系统越大，其中局部组件失效的概率就越大。在长时间运行期间，失效，修复，再失效可以看成是一个反复的过程。在一个包含成千上万个节点的系统中，我们几乎总是可以假定某些东西发生了失效<sup>[7]</sup>。此时，如果采用的是简单的停止-修复错误处理策略，对于这样一个庞大的集群系统，最终将花费大量时间在错误恢复上而不是正常的任务执行<sup>[8]</sup>。
- 如果系统可以容忍某些失败的节点，而使整体继续工作，则对系统运维帮助极大。例如，支持滚动升级（参阅第4章）：每次重启一个节点，而集群总体对外不中断服务。在云环境中，如果发现某台虚拟机有问题，可以将其关闭然后重新启动另一个。
- 对于全球分散部署的多数据中心（使用户访问地理靠近的数据中心，从而降低延迟），通信很可能经由广域网，与本地网络相比，速度更慢且更加不可靠。而高性能计算通常假设所有节点位置靠近、紧密连接。

要使分布式系统可靠工作，就必然面临部分失效，这就需要依靠软件系统来提供容错机制。换句话说，我们需要在不可靠的组件之上构建可靠的系统。另外，正如在第1章“可靠性”所讨论的，世界上不存在完美的可靠性，我们需要的现实可行的保证。

即使对于只有几个节点的小型系统，也很有必要审视部分失效问题。在一个小系统中，很可能大部分组件在大部分时间都正常工作，但迟早某一天有一部分系统会出现故障，此时软件必须可以有效处理。这里要强调的是，故障处理是软件设计的重要组成部分。作为系统运维者，需要知道在发生故障时，系统的预期行为是什么。

不能假定故障不可能发生而总是期待理想情况。最好仔细考虑各种可能的出错情况，包括那些小概率故障，然后尝试人为构造这种测试场景来充分检测系统行为。可以说，在分布式系统中，怀疑，悲观和偏执狂才能生存。

## 不可靠的网络

正如第二部分所介绍的，本书关注的主要是分布式无共享系统，即通过网络连接的多个节点。网络是跨节点通信的唯一途径，我们还假定每台机器都有自己的内存和磁盘，一台机器不能直接访问另一台机器的内存或磁盘除非通过网络向对方发出请求。



## 基于不可靠的组件构建可靠的系统

直观来看，系统的可靠性应该取决于最不可靠的组件（即最薄弱的环节）。然而，对于计算机系统来讲，事实并非如此，在低可靠性部件上构建高可靠的系统一直是计算机领域的惯用手段<sup>[11]</sup>。例如：

- 纠错码可以在各种通信链路上准确传输数据，包括那些可能偶尔传输出错的情况，例如无线网络出现信号干扰<sup>[12]</sup>。
- IP（Internet协议）层本身并不可靠，可能会出现丢包、延迟、重复发送以及乱序等情况。但 TCP（传输控制协议）在IP之上提供了更可靠的传输层，可以保证丢失的数据包被重传，消除重复包，包的顺序以发送的顺序重新组合等。

需要指出，系统整体虽然变得比底层组件更为可靠，但可靠性总有其现实上限。例如，纠错码只能处理某几个位错误，但如果信号被彻底干扰，就很难保证最终可以正确接受多少数据<sup>[13]</sup>。而TCP虽然能够提供重传、重组等，但是它肯定不能神奇地消除网络传输延迟。

尽管系统并不完美，但这样的高可靠系统仍然非常有用，它可以帮我们处理底层一些棘手的故障，使其他故障更容易理解和进一步处理。我们将在第12章进一步探讨这个问题。

首先要说明，无共享并不是构建集群系统的唯一方式，但它却是构建互联网服务的主流方式。主要是由于以下几个原因：由于不需要专门的硬件因此成本相对低廉，可以采用通用的商品化硬件，可以采用跨区域的多数据中心来实现高可靠性。

互联网以及大多数数据中心的内部网络（通常是以太网）都是异步网络。在这种网络中，一个节点可以发送消息（数据包）到另一个节点，但是网络并不保证它什么时候到达，甚至它是否一定到达。发送之后等待响应过程中，有很多事情可能会出错（见图8-1所示的例子）：

1. 请求可能已经丢失（比如有人拔掉了网线）。
2. 请求可能正在某个队列中等待，无法马上发送（也许网络或接收方已经超负荷）。
3. 远程接收节点可能已经失效（例如崩溃或关机）。



4. 远程接收节点可能暂时无法响应（例如正在运行长时间的垃圾回收，请参阅本章后面的“进程暂停”）。
5. 远程接收节点已经完成了请求处理，但回复却在网络中丢失（例如网络交换机配置错误）。
6. 远程接收节点已经完成了请求处理，但回复却被延迟处理（例如网络或者发送者的机器超出负荷）。

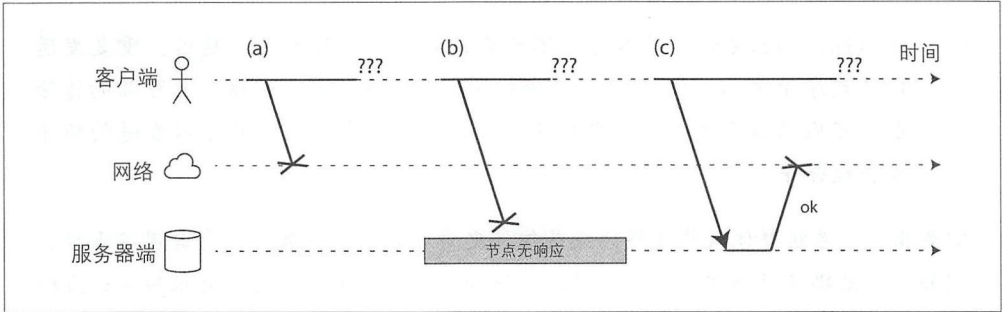


图8-1：如果请求没有得到响应，无法区分：（a）请求丢失；（b）远程节点关闭；（c）响应丢失

发送者甚至不清楚数据包是否完成了发送，只能选择让接收者来回复响应消息，但回复也有可能丢失或延迟。这些问题在一个异步网络中无法明确区分，发送者拥有的唯一信息是，尚未收到响应，但却无法判定具体原因。

处理这个问题通常采用超时机制：在等待一段时间之后，如果仍然没有收到回复则选择放弃，并且认为响应不会到达。但是，即使判定超时，仍然并不清楚远程节点是否收到了请求（一种情况，请求仍然在某个地方排队，即使发送者放弃了，但最终请求会发送到接收者）。

## 现实中的网络故障

我们已经有几十年的计算机网络构建经验，人们也许想当然地会认为现在已经掌握了足够的技能使网络变得非常可靠。但是，情况似乎并非如此。

一些系统研究和大量的侧面证据表明，网络问题出人意料地普遍，包括那些由公司运营的数据中心<sup>[14]</sup>。一家中型数据中心完成的调查发现，每月大约有12次网络故障，其中有一半涉及单台机器，有一半甚至是整个机架断网<sup>[15]</sup>。另一项研究分析了机架式交换机、汇聚层交换机和负载均衡器等组件的故障率<sup>[16]</sup>，结果发现增加冗余网络设备并



不会像期望的那样可以显著减少故障，主要是因为无法有效防范人为错误（例如配置错误），而后者是造成网络中断的主要原因。

公共云服务如AWS 因多次出现临时性网络故障而影响颇大<sup>[14]</sup>，管理良好的私有数据中心网络或许更为稳定一些。尽管如此，所有人都曾遭受网络问题的困扰：例如，交换机软件升级可能会触发网络拓扑重新配置，在此期间数据包的延迟显著增加，甚至超过了一分钟<sup>[17]</sup>；有报道发现鲨鱼咬破了海底电缆<sup>[18]</sup>。其他令人惊讶的故障还包括网络接口会丢弃所有入向数据包，但可以成功发送出向数据包<sup>[19]</sup>，原因仍然在于网络接口配置问题。



### 网络分区

当网络的一部分由于网络故障而与其他部分断开，称之为网络分区或网络分割。为避免与第6章存储系统的分区（分片）产生混淆，本书采用更通用的术语“网络故障”。

即使网络故障在你的环境中比较少见，但故障可能发生也要求软件需要能够处理它们。事实上，只要有网络通信，就可能会出现故障，这一点始终无法彻底避免。

如果没有处理或者测试网络故障，可能会发生意想不到的后果。例如，集群可能会死锁，即使网络恢复了也无法提供服务<sup>[20]</sup>，甚至可能误删除数据<sup>[21]</sup>。如果触发了一些软件未定义的情形，则发生任何意外都不奇怪。

处理网络故障并不意味着总是需要复杂的容错措施：假定你的网络通常非常可靠，而万一出现问题，一种简单的方法是对用户提示错误信息。但前提是，必须非常清楚接下来软件会如何应对，以确保系统最终可以恢复。我们推荐有计划地人为触发网络问题，目的是测试系统的反应情况（这也是Chaos Monkey系统背后的想法，请参阅第1章“可靠性”）。

## 检测故障

许多系统都需要自动检测节点失效这样的功能。例如：

- 负载均衡器需要避免向已失效的节点继续分发请求（即将其下线处理）。
- 对于主从复制的分布式数据库，如果主节点失败，需要将某个从节点提升为主节点（参阅第5章“处理节点失效”）。不过，由于网络的不确定性很难准确判断节点是否确实失效。



- 然而不幸的是，由于网络的不确定性使得判断节点是否失效非常困难；而只有在某些特定场景下，或许你可以明确知道哪里出错了。
- 假设可以登录节点，但发现服务进程没有侦听目标端口（例如，由于进程已经崩溃），那么操作系统会返回RST或FIN标志的数据包来辅助关闭或拒绝TCP连接。但是，如果节点在处理请求的过程中发生了崩溃，则很难知道该节点实际处理了多少数据<sup>[22]</sup>。
- 如果服务进程崩溃（或被管理员杀死），但操作系统仍正常运行，可以通过脚本通知其他节点，以便新节点来快速接管而跳过等待超时。HBase采用了这样的方法<sup>[23]</sup>。
- 如果有权访问数据中心网络交换机，则可以通过管理接口查询是否存在硬件级别的链路故障（例如远程节点掉电）。不过，该方法也有局限性，例如通过互联网连接，或者是处于共享数据中心而没有访问交换机的权限，以及由于网络问题而根本无法登录管理界面等。
- 如果路由器已经确认目标节点不可访问，则会返回ICMP“目标不可达”数据包来回复请求。但是，路由器本身并不具有什么神奇的检测能力，从这一点来讲，它和网络上其他节点并无本质区别。

能快速告之远程节点的关闭状态自然有用，但也不是万能的。例如，即使TCP确认一个数据包已经发送到目标节点，但应用程序也可能在处理完成之前发生崩溃。如果你想知道一个请求是否执行成功，就需要应用级别的回复<sup>[24]</sup>。

总之，如果出现了问题，你可能会在应用堆栈的某个级别拿到了一个关于错误的回复，但最好假定最终收不到任何错误报告。接下来尝试重试（TCP重试是透明的，但也可以在应用级别重试），等待超时之后，如果还是没有收到响应，则最终声明节点已经失效。

## 超时与无限期的延迟

如果超时是故障检测唯一可行的方法，那么超时应该设多长呢？不幸的是没有标准的答案。

设置较长的超时值意味着更长时间的等待，才能宣告节点失效（在此期间，用户只能等待或者拿到错误信息）。较短的超时设置可以帮助快速检测故障，但可能会出现误判，例如实际上节点只是出现暂时的性能波动（由于节点或网络上的高负载峰值），结果却被错误地宣布为失效。



如果节点实际上是活着（例如，正在发送一封电子邮件），过早地将其声明为失效会带来一些问题，新节点会尝试接管，然后出现一些操作在两个节点上执行两次。我们将在后面详细讨论这个问题。

当一个节点被宣告为失效，其承担的职责要交给到其他节点，这个过程会给其他节点以及网络带来额外负担，特别是如果此时系统已经处于高负荷状态。例如节点只是负载过高而出现了响应缓慢，转移负载到其他节点可能会导致失效扩散，产生级联扩大效应，在极端情况下，所有节点都宣告对方死亡，造成服务处于事实停止状态。

设想一个虚拟的系统，其网络可以保证数据包的最大延迟在一定范围内：要么在时间 $d$ 内完成交付，要么丢失。此外，假定一个非故障节点总能够在一段时间 $r$ 内完成请求处理。此时，可以确定成功的请求总能够在 $2d+r$ 时间内收到响应，如果在此时间内没有收到响应，则可以推断网络或者远程节点发生了失效，那么 $2d+r$ 是一个理想的超时设置。

然而事实上绝大多数系统都没有类似的保证：异步网络理论上的延迟无限大（即使尽力发送数据包，但数据包到达时间并没有上确界），多数服务端也无法保证在给定的某个时间内一定完成请求处理（参阅本章后面的“响应时间保证”）。如果超时设置太小，只需要一个短暂的网络延迟尖峰就会导致包超时进而将系统标记为失效。

## 网络拥塞与排队

驾车时往往由于交通堵塞，导致行车时间变化很大。同样，计算机网络上数据包延的变化根源往往在于排队<sup>[25]</sup>：

- 当多个不同节点同时发送数据包到相同的目标节点时，网络交换机会出现排队，然后依次将数据包转发到目标网络（见图8-2）。如果网络负载过重，数据包可能必须等待一段时间才能获得发送机会（即网络拥塞）。如果数据量太大，交换机队列塞满，之后的数据包则会被丢弃，网络还在运转，但会引发大量数据包重传。
- 当数据包到达目标机器后，如果所有CPU核都处于繁忙状态，则网络数据包请求会被操作系统排队，直到应用程序能够处理。根据机器的配置和负载情况，这里也会引入一段不确定的等待时间。
- 在虚拟化环境下，CPU核会切换虚拟机，从而导致正在运行的操作系统会突然暂停几十毫秒。在这段时间，客户虚机无法从网络中接收任何数据，入向的包会被虚拟机管理器<sup>[26]</sup>排队缓冲，进一步增加了网络延迟的不确定性。
- TCP执行流量控制（也称为拥塞消除，或背压）时，节点会主动限制自己的发送





速率以避免加重网络链路或接收节点负载<sup>[27]</sup>。这意味着数据甚至在进入网络之前，已经在发送方开始了排队。

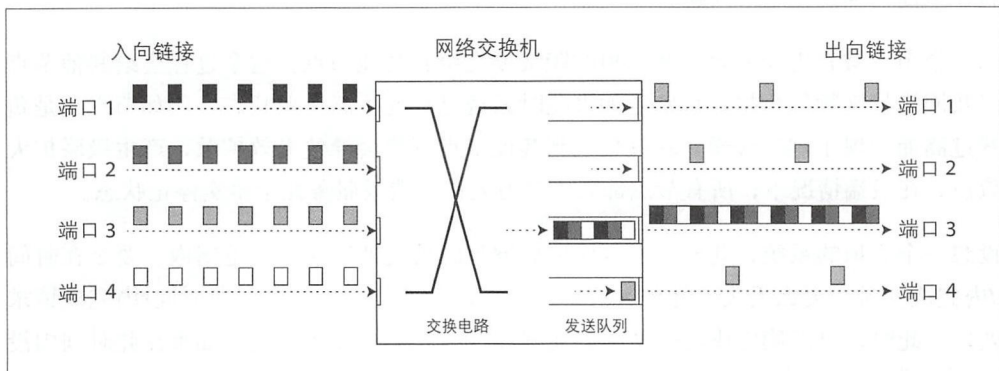


图8-2：多台节点同时发送流量到同一目标，交换机队列可能被填满。图中端口1、2和4都在发送数据包到端口3

而且，对于TCP如果在某个超时范围内（根据往返时间来推算）没有收到确认，则认为数据包已经丢失进而触发自动重传。虽然这一过程对应用程序透明，但肯定会引入额外的延迟（等待超时到期，等待重传的数据包得到确认）。

## TCP与UDP

一些对延迟敏感的应用程序（如视频会议和IP语音VoIP）使用了UDP而不是TCP。这是一个可靠性与可变性之间的权衡考虑。UDP不支持流量控制也不会重传丢失的数据包，因此可以避免一些网络延迟不确定的因素（但它仍受队列和调度延迟等方面的影响）。

如果延迟或丢弃的数据价值不大，UDP是个不错的选择。例如，IP电话情况可能没有足够的时间重传丢失的数据包，或者重传并没有太大的意义，应用程序会采用静音填充丢失的位置（出现短暂声音中断），然后数据流必须尽快向前继续，由人为即通话双方来沟通重试。类似这样，“请你再说一遍好吗？刚刚声音没有听到。”

所有以上因素都会造成网络延迟的变化或者不确定性。当系统还有足够的处理能力，排队之后可以快速处理；但当系统接近其最大设计上限时，系统负载过高，队列深度就会显著加大，排队对延迟的影响变得特别明显。

在公有云和多租户数据中心中，许多客户共享网络资源包括交换机、机器的网卡以及

CPU（如虚拟机）等。批处理工作负载例如MapReduce（参阅第10章）很容易使网络带宽达到饱和。通常用户无法控制或者没有权限观测其他用户对共享资源的使用情况，如果不巧旁边有个疯狂的邻居（如虚拟机）正在大量使用资源，网络延迟就会波动很大<sup>[28,29]</sup>。

在这种环境下，只能通过实验方式一步步设置超时。先在多台机器上，多次测量网络往返时间，以确定延迟的大概范围；然后结合应用特点，在故障检测与过早超时风险之间选择一个合适的中间值。

更好的做法是，超时设置并不是一个不变的常量，而是持续测量响应时间及其变化（抖动），然后根据最新的响应时间分布来自动调整。可以用Phi Accrual故障检测器<sup>[30]</sup>完成，该检测器目前已在Akka和Cassandra<sup>[31]</sup>中使用。TCP的重传超时也采用了类似的机制<sup>[27]</sup>。

## 同步与异步网络

如果网络层可以在规定的延迟内保证完成数据包的发送，且不会丢弃数据包，那么分布式系统就会简单很多。为什么我们不能考虑在硬件层解决这个问题呢？使网络足够可靠，然后软件就无需如此担心。

为了回答这个问题，可以将数据中心网络与传统的固定电话网络（非移动蜂窝，非VoIP）进行对比分析，首先后者非常可靠：语音延迟和掉话现象极为罕见。这样的固定电话网络需要有持续端到端的低延迟和足够的带宽来传输音频数据。计算机网络能否实现类似的高可靠性和确定性呢？

当通过电话网络拨打电话时，系统会动态建立一条电路：在整个线路上为呼叫分配一个固定的、带宽有保证通信链路，该电路一直维持到通话结束<sup>[32]</sup>。例如，ISDN网络固定以每秒4000帧的速率运行。当新的呼叫建立后，在每个帧（每个方向）内分配16bit的空间，在通话期间，每一方都可以保证在250μs内完成发送一条16bit的音频数据<sup>[33,34]</sup>。

这种网络本质是同步的：即使数据中间经过了多个路由器，16bit空间在电路建立时已经在网络中得到预留，不会受到排队的影响。由于没有排队，网络最大的端到端延迟是固定的。我们称之为有界延迟。

### 网络延迟是否可预测？

请注意，固定电话网络中的电路与TCP连接存在很大不同：电路方式总是预留固定带

宽，在电路建立之后其他人无法使用；而TCP连接的数据包则会尝试使用所有可用的网络带宽。TCP可以传送任意大小可变的数据块（例如电子邮件或网页），它会尽力在最短的时间内完成数据发送。而当TCP连接空闲时，通常不占用任何带宽<sup>注2</sup>。

假设数据中心网络和互联网基于电路交换网络，一旦电路建立完成则可以保证最大往返时间。然而，事实是以太网和IP都是基于分组交换协议，这种协议注定受到排队的影响，从而导致网络延迟不确定，在这些协议里完全没有电路的概念。

那为什么数据中心网络和互联网采用分组交换呢？答案是，它们针对突发流量进行了很多优化。电路非常适合音频或视频通话，通话期间只需每秒传送固定数量的数据。但对于访问网页，发送电子邮件或传输文件等无法事先确定带宽需求，我们只是希望它尽快完成。

如果你想通过电路链接来传输文件，将不得不预估一个待分配的带宽。如果预估值太低，传输速度就特别缓慢，甚至无法实际可用；如果预估带宽太高，电路甚至无法完成建立（因为如果无法预留所需的带宽，电路就无法建立）。所以，对于突发数据的传输，电路网络无法充分利用网络容量，导致发送缓慢。相比之下，TCP动态调整传输速率则可以充分利用所有可用的网络容量。

曾经也有一些尝试建立支持电路交换和分组交换的混合网络，比如ATM<sup>注3</sup>。InfiniBand网络有一些相似之处<sup>[35]</sup>：它在链路层实现端到端的流量控制，从而减少了网络中的排队，但仍可能会由于链路拥塞而影响延迟<sup>[36]</sup>；然后通过服务质量（QoS，数据包的优先级和调度）和准入控制（限制发送速率），最终可以在分组网络上模拟电路交换，或者说提供统计意义上的有限延迟<sup>[25,32]</sup>。

但是，目前此类QoS在多租户数据中心、公有云和广域网<sup>注4</sup>中并未启用。总之，当前广泛部署的技术无法为我们提供延迟或可靠性方面的硬件级保证，我们必须假设会出现网络拥塞，排队和无上限的延迟。基于此，超时设置并不存在某个绝对正确的值，而是需要通过实验的方式来确定。

---

注2： 当TCP keep alive功能启用时，由少量的keep alive数据包。

注3： 异步传输模式（Asynchronous Transfer Mode，ATM）在20世纪80年代曾是以太网的竞争对手<sup>[32]</sup>，但最终除了电话网核心交换机之外ATM并没有大规模部署。这个缩略语与自动收贷机（或者自动取款机）相同，但二者并无任何关系。或许，在某个平行的世界里，互联网是以ATM为基础构建的，互联网视频通话更可靠，也会遭遇丢包和延迟困扰。

注4： 互联网服务提供商之间的对等协议，以及通过边界网关协议（BGP）所建立路由的协议与电路交换更为相似。在这个级别上，可以购买专用带宽。然而，互联网路由与主机直接的网络连接并不在相同的工作层面，其范围、尺度更广。



## 延迟与资源利用率

从更广泛的意义讲，可以将延迟的波动归为动态资源分区的结果。

假设两台电话交换机之间有一根线路可同时支撑10 000个呼叫。线路上的每条电路都占用其中一个槽位，因此，可以将线路视为最多允许10 000个并发用户所共享的资源。只是资源是以静态方式分配的：即使现在你是线路上唯一的通话者，且剩余的其他9 999个槽位都是空闲状态，你的电路也只能使用固定数量的带宽。

相反，互联网则是动态分享网络带宽。多个发送方互相竞争，以尽快地通过网络发生他们的数据，网络交换机也是动态决定该发送哪个数据包。这种方法存在排队缺点，但优点是最大限度地利用了带宽。线路本身对应一个固定的成本，如果可以更充分的使用，那么发送每个字节的平均摊薄成本自然就会下降。

CPU也是类似的情况：如果多个线程动态共享CPU核，当有线程正在CPU上运行时，其他线程必须在操作系统的运行队列上等待，这样会出现长短不一的执行暂停。但是，相比于为每个线程静态分配固定的CPU执行周期（参见本章后面的“响应时间保证”），它可以更好地利用硬件，这是虚拟化技术非常重要的动机。

如果资源总是静态分配（例如，专用硬件和预留带宽分配），则某些环境下可以保证延迟的确定性。但是，这是以降低资源使用率为代价的，换句话说，其成本过于昂贵。而多租户、动态资源分配方式则可以提供更高的资源使用率，因而成本更低，当然也引入了可变延迟的缺点。

简言之，网络中的可变延迟并不是一种自然规律，只是成本与收益相互博弈的结果。

## 不可靠的时钟

时钟和计时非常重要。有许多应用程序以各种方式依赖于时钟，例如：

1. 某个请求是否超时了？
2. 某项服务的99%的响应时间是多少？
3. 在过去的五分钟内，服务平均每秒处理多少个查询？
4. 用户在我们的网站上浏览花了多长时间？
5. 这篇文章什么时候发表？

6. 在什么时间发送提醒邮件?
7. 这个缓存条目何时过期?
8. 日志文件中错误消息的时间戳是多少?

上述1~4测量的持续时间（例如请求发送与响应接收之间的时间间隔），而5~8所描述的某个时间点（在特定日期，特定时间发生的事件）。

在分布式系统中，时间总是件棘手的问题，由于跨节点通信不可能即时完成，消息经由网络从一台机器到另一台机器总是需要花费时间。收到消息的时间应该晚于发送的时间，但是由于网络的不确定延迟，精确测量面临着很多挑战。这些情况使得多节点通信时很难确定事情发生的先后顺序。

而且，网络上的每台机器都有自己的时钟硬件设备，通常是石英晶体振荡器。这些设备并非绝对准确，即每台机器都维护自己本地的时间版本，可能比其他机器稍快或更慢。可以在一定程度上同步机器之间的时钟，最常用的方法是网络时间协议（Network Time Protocol, NTP），它可以根据一组专门的时间服务器来调整本地时间<sup>[37]</sup>，时间服务器则从精确更高的时间源（如GPS接收机）获取高精度时间。

## 单调时钟与墙上时钟

现代计算机内部至少有两种不同的时钟：一个是墙上时钟（或称钟表时间），一个是单调时钟。虽然它们都可以衡量时间，但要仔细区分二者，本质上他们是服务于不同的目的。

### 墙上时钟

墙上时钟根据某个日历（也称为墙上时间）返回当前的日期与时间。例如，Linux的`clock_gettime(CLOCK_REALTIME)`<sup>注5</sup>和Java中的`System.currentTimeMillis()`会返回自纪元1970年1月1日（UTC）以来的秒数和毫秒数，不含闰秒。而有些系统则使用其他日期作为参考点。

墙上时钟可以与NTP同步。但是，如下一节所述，这里还存在一些奇怪问题。特别是，如果本地时钟远远快于NTP服务器，强行重置之后会跳回到先前的某个时间点。这种跳跃以及经常忽略闰秒，导致其不太适合测量时间间隔<sup>[38]</sup>。

---

注5： 尽管这里有所谓的实时，但其实与实时操作系统无关，请参阅本章后面的“响应时间保证”。

曾经墙上时钟粗糙的精度也被人诟病，例如在较早的Windows系统上只能提供10毫秒的精度<sup>[39]</sup>。不过在较新的系统上，精度已经不成问题。

## 单调时钟

单调时钟更适合测量持续时间段（时间间隔），例如超时或服务的响应时间：Linux上的`clock_gettime(CLOCK_MONOTONIC)`和Java中的`System.nanoTime()`返回的即是单调时钟。单调时钟的名字来源于它们保证总是向前（而不会出现墙上时钟的回拨现象）。

可以在一个时间点读取单调时钟的值，完成某项工作，然后再次检查时钟。时钟值之间的差值即两次检查之间的时间间隔。注意，单调时钟的绝对值并没有任何意义，它可能电脑启动以后经历的纳秒数或者其他含义。因此比较不同节点上的单调时钟值毫无意义，它们没有任何相同的基准。

如果服务器有多路CPU，则每个CPU可能有单独的计时器，且不与其他CPU进行同步。由于应用程序的线程可能会调度到不同的CPU上，此时，操作系统会补偿多个计时器之间的偏差，从而为应用层提供统一的单调递增计时。不过最好还是对这种偏差补偿持谨慎态度<sup>[40]</sup>。

如果NTP检测到本地石英比时间服务器上更快或者更慢，NTP会调整本地石英的震动频率（这被称为摆动）。默认情况下，NTP允许速率加快或减慢的最大幅度为0.05%，但NTP并不会直接调整单调时钟向前或者回拨。单调时钟的精度通常很高，在如今大多数系统中，可以测量几微秒或更短的时间间隔。

在分布式系统中，可以采用单调时钟测量一段任务的持续时间（例如超时），它不假定节点间有任何的时钟同步，且可以容忍轻微测量误差。

## 时钟同步与准确性

单调时钟不需要同步，但是墙上时钟需要根据NTP服务器或其他外部时间源做必要的调整。然而，我们获取时钟的方法并非预想那样可靠或准确，硬件时钟和NTP可能会出现一些莫名其妙的现象。举几个例子：

- 计算机中的石英钟不够精确，存在漂移现象（运行速度会加快或减慢）。时钟漂移主要取决于机器的温度。谷歌假设其服务器的时钟偏移为200 ppm（百万分之一）<sup>[41]</sup>，相当于如果每30秒与服务器重新同步一次，则可能出现的最大偏差为



6毫秒，或者每天一次同步，则最大偏差为17秒。即使其他一切工作正常，漂移问题也限制了可以达到的最佳精度。

- 如果时钟与NTP服务器的时钟差别太大，可能会出现拒绝同步，或者本地时钟将被强制重置<sup>[37]</sup>。在重置前后应用程序观察可能会看到时间突然倒退或突然跳跃的现象。
- 由于某些原因，如果与NTP服务器链接失败（如防火墙），可能会很长一段时间没有留意到错误配置最终导致同步失败。有证据表明，实践中确实发生这类问题。
- NTP同步会受限于当时的网络环境特别是延迟，如果网络拥塞、数据包延迟变化不定，则NTP同步的准确性会受影响。实验表明，当通过互联网进行同步时，可能会产生至少35毫秒的偏差，最坏时（如网络抖动出现包发送失败）则可能超过1秒。取决于具体的参数配置情况，如果网络延迟特别严重，则NTP客户端可能会被迫放弃同步。
- 一些NTP服务器本身出现故障、或者配置错误，其报告的时间可能存在数小时的偏差<sup>[43,44]</sup>。NTP客户端往往比较稳定，可以同时查询多个服务器并忽略掉异常值。尽管如此，相信互联网上一个陌生的服务器告诉你的时间，还是值得警惕。
- 闰秒会产生一分钟为59秒或61秒的现象，这会在使一些对闰秒毫无防范的系统出现混乱<sup>[45]</sup>。闰秒曾经使许多大型系统崩溃，过去的教训表明，不经意间许多系统已经埋下了对时钟不正确处理的隐患。而处理闰秒的推荐方式是，不管NTP服务器具体如何实现，在NTP服务器汇报时间时故意做些调整，目的是在一天的周期内逐步调整闰秒（称为拖尾）<sup>[47-49]</sup>。
- 在虚拟机中，由于硬件时钟也是被虚拟化的，这对于需要精确计时的应用程序提出了额外的挑战<sup>[50]</sup>。当虚拟机共享一个CPU核时，每个虚拟机会出现数十毫秒内的暂停以便切换客户虚拟机。但从应用的角度来看，这种停顿会表现为时钟突然向前发生了跳跃<sup>[26]</sup>。
- 如果运行在未完全可控的设备上（例如，移动设备或嵌入式设备），需要留意不能完全相信设备上的硬件时钟。某些用户会故意将其硬件时钟设置为错误的日期和时间，例如为了规避游戏的时间限制。因此，获取的时钟可能是某个过去的、或者将来的时间。

如果确实需要投入大量资源，是可以达到非常高的时钟精度。例如，针对金融机构的欧洲法规草案MiFID II就明确要求所有高频交易基金必须在UTC时间100微秒内同步时钟，以便调试“崩盘”等市场异常并检测市场操纵等违规行为<sup>[51]</sup>。



高精度的时钟可以采用GPS接收机，精确时间协议（PTP）<sup>[52]</sup>并辅以细致的部署和监测。但通常也意味着大量的资源投入和技术门槛，并持续监控时钟同步可能出现的错误情况。例如NTP守护程序配置错误，防火墙阻止NTP通信等，避免导致时钟误差迅速变大。

## 依赖同步的时钟

时钟虽然看起来简单，但却有不少使用上的陷阱：一天可能不总是86 400秒，时钟会向后回拨，一个节点上的时间可能与另一个节点上的时间完全不同。

本章前面讨论了网络丢包和数据包延迟。对此，我们的建议是，即使网络在大多数情况下表现良好，软件也必须假设网络偶尔会出现故障，因此要有对应的措施来妥善处理这些故障。时钟也是如此：尽管大多数时间工作很好，但仍需以备不测。

假如一台机器的CPU出现了故障或者网络有问题，系统可能根本无法工作，所以很快就会被注意到进而得到修复，但时钟问题却不那么容易被及时发现。如果石英时钟有缺陷，或者NTP客户端配置错误，最后出现了时间偏差，对大多数功能可能并没太大影响。但对于一些高度依赖于精确时钟的软件，出现的后果可能是隐式的，或许会丢失一小部分数据而不是突然的崩溃<sup>[53,54]</sup>。

因此，如果应用需要精确同步的时钟，最好仔细监控所有节点上的时钟偏差。如果某个节点的时钟漂移超出上限，应将其宣告为失效，并从集群中移除。这样的监控的目的是确保在造成重大影响之前尽早发现并处理问题。

## 时间戳与事件顺序

对于一个常见的功能：跨节点的事件排序，如果它高度依赖时钟计时，就存在一定的技术风险。例如，两个客户端同时写入分布式数据库，谁先到达？哪一个操作是最新的呢？

图8-3给出了这样的危险例子，即多主节点复制的分布式数据库高度依赖于墙上时钟（图5-9是另一个类似的例子）。客户端A在节点1上写入 $x = 1$ ，写入被复制到节点3；客户端B在节点3上增加 $x$ （现在 $x = 2$ ）。最后，这两个写入都被复制到节点2上。

在图8-3中，写入被复制到其他节点时，会根据源写入节点上的墙上时钟来标记时间戳。在该例子中，时钟同步机制稳定工作，节点1和节点3之间的时钟偏差小于3ms，这比实践中的多数情况可能都要更好。



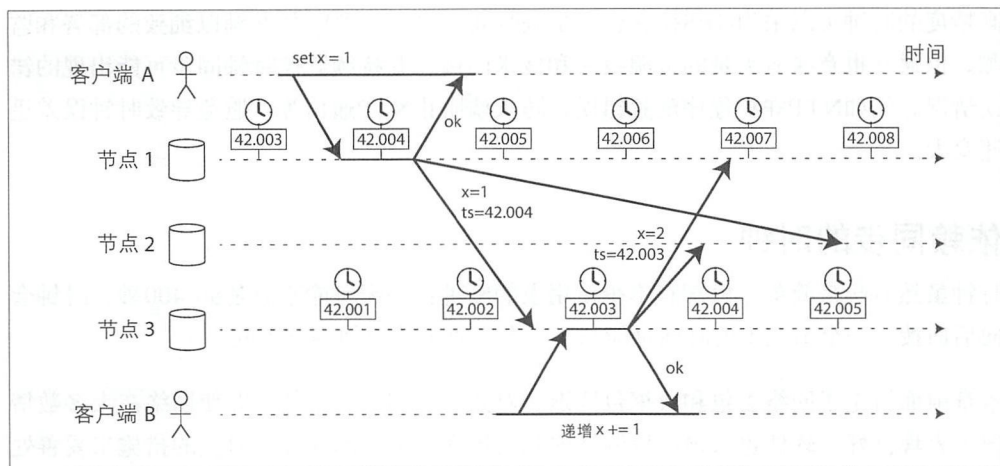


图8-3：客户端B的写入比客户端A的写入要晚，但是B写入的时间戳却更早

但是，这样的时间戳却不能正确排序事件：写入  $x = 1$  的时间戳为42.004秒，写入  $x = 2$  虽然后续发生，时间戳却是42.003s。当节点2收到这两个事件时，会根据时间戳错误地判断  $x = 1$  是最新值，然后决定丢弃  $x = 2$ ，这就导致客户端B的增量操作丢失。

这种冲突解决策略被称为最后写入获胜（LWW），在多主节点以及无主节点复制数据库（如Cassandra<sup>[53]</sup>和Riak<sup>[54]</sup>）中广泛使用（参阅第5章“最后写入获胜”）。有些实现会在客户端生成时间戳而非服务器端，但无论如何没有改变LWW的根本问题：

- 数据库写入可能会奇怪地丢失：明明后续发生的写操作却没法覆盖另一个较早的值，原因是后者节点的时钟太快了<sup>[54,55]</sup>。这会导致一些数量未知的数据被悄悄地丢弃，并且不会向应用报告任何错误。
- LWW无法区分连续快速发生的连续写操作（图8-3中客户端A写入之后才发生了客户端B的增量操作）和并发写入（每个写操作都不依赖于其他写）。需要额外的因果关系跟踪机制（例如版本向量）来防止因果冲突（参阅第5章“检测并发写”）。
- 由于时钟精度的限制（例如毫秒级），两个节点可能各自独立产生了完全相同的时间戳。为了解决这样的冲突，需要一个额外的仲裁值（可以简单地引入一个大的随机数），但该方法还是无法区分因果关系<sup>[53]</sup>。

因此，通过保持“最新”值并丢弃其他值来解决冲突看似不错，但要注意，“最新”的定义如果取决于墙上时钟就会引入偏差。即使采用了NTP同步时钟，依然可能会出





现在时间戳100毫秒时（根据发送者的时钟）发送了某个数据包，却在时间戳99毫秒（根据接收者的时钟）到达，这看起来好像是数据包还没发送就先到达了。

那么NTP时钟同步能否做到极高的精度从而避免这种错误的顺序问题呢？很难做到。因为除了石英漂移等误差来源之外，NTP同步精度本身要受限于所在的网络延迟。要达到正确的排序，需要时钟源精度要远远高于被测量的对象（即网络延迟）。

对于排序来讲，基于递增计数器而不是振荡石英晶体的逻辑时钟<sup>[56,57]</sup>是更可靠的方式（参见第5章“检测并发送”）。逻辑时钟并不测量一天的某个时间点或时间间隔，而是事件的相对顺序（事件发生的相对前后关系）。与之对应的，墙上时钟和单调时钟都属于物理时钟。我们会在第9章“顺序保证”中继续讨论顺序问题。

## 时钟的置信区间

或许墙上时钟会返回微秒甚至纳秒级别的信息，但是这种精度的测量值其实并不可信。如前所述，石英漂移问题可导致偏差高达几毫秒，即使每分钟都与本地网络的NTP服务器进行同步，也无法保证上述精度。如果使用公共互联网上的NTP服务器，最好的精度也只能到几十毫秒，而一旦出现网络拥塞，偏差很容易就超过100毫秒<sup>[57]</sup>。

因此，我们不应该将时钟读数视为一个精确的时间点，而更应该视为带有置信区间的时间范围。例如，系统可能有95%的置信度认为目前时间介于10.3 ~ 10.5秒之间<sup>[58]</sup>。如果我们可完全相信的精度为 $\pm 100$ 毫秒，那么时间戳中那些微秒级的读数并无实际意义。

可以根据具体的时间源来推算出时钟误差的上限。如果节点上直接装有GPS接收器或原子（铯）时钟，那它的误差范围通常可查询制造商的手册。如果节点是从服务器获取时间，则不确定性取决于上次服务器同步以来的石英漂移范围，加上NTP服务器的不确定性，再加上与服务器之间的网络往返时间（对于第一次同步，我们假定服务器完全可信）。

可惜大多数系统并不提供这种误差查询接口。例如，当调用`clock_gettime()`时，返回值没有任何误差信息，所以无法确切知道置信区间应该是五毫秒还是五年。

这里有趣的是Google Spanner中的`TrueTime` API<sup>[41]</sup>，它会明确地报告本地时钟的置信区间。当查询当前时间时，你会得到两个值：[不早于，不晚于]分别代表误差的最大偏差范围。基于上述两个时间戳，可以知道实际时间应在其范围之内。该间隔的范围主要取决于本地石英钟最后与高精时钟源同步后所经历的时间长短。



## 全局快照的同步时钟

在第7章“快照隔离与可重复读”中，我们介绍了快照隔离，它广泛用于小数据量、快速读写的事务以及大数据量，长时间运行的只读事务（例如备份或分析），可以在数据库的某个一致状态上不需加锁、不违背读写隔离性的前提下高效支持只读事务。

常见的快照隔离实现中需要单调递增事务ID。如果写入发生在快照之后（即写入具有比快照更大的事务ID），那么该写入对于快照不可见。在单节点数据库上，一个简单的计数器足以生成事务ID。

但是，当数据库分布在多台机器上（可能跨越多个数据中心）时，由于需要复杂的协调以产生全局的、单调递增的事务ID（跨所有分区）。事务ID要求必须反映因果关系：事务B如果要读取事务A写入的值，则B的事务ID必须大于A的事务ID，否则快照将不一致。考虑到大量、频繁的小数据包，在分布式系统中创建事务ID通常会引入瓶颈<sup>注6</sup>。

能否使用同步后的墙上时钟作为事务ID呢？如果我们能够获得足够可靠的同步时钟，自然它可以符合事务ID属性要求：后发生的事务具有更大的时间戳。然而问题还是时钟精度的不确定性。

Google Spanner采用以下思路来实现跨数据中心的快照隔离<sup>[59,60]</sup>。它根据TrueTime API返回的时钟置信区间，并基于以下观察结果：如果有两个置信区间，每个置信区间都包含最早和最新可能的时间戳（ $A = [A_{\text{earliest}}, A_{\text{latest}}]$ 和 $B = [B_{\text{earliest}}, B_{\text{latest}}]$ ），且这两个区间没有重叠（即 $A_{\text{earliest}} < A_{\text{latest}} < B_{\text{earliest}} < B_{\text{latest}}$ ），那么可以断定B一定发生在A之后。只有发生了重叠，A和B发生顺序才无法明确。

为了确保事务时间戳反映因果关系，Spanner在提交读写事务之前故意等待置信区间的长度。这样做的目的是，确保所有读事务要足够晚才发生，避免与先前的事务的置信区间产生重叠。而为了尽量缩短潜在的等待时间，Spanner需要使时钟的误差范围尽可能的小，为此，Google在每个数据中心都部署了一个GPS接收器或原子钟，保证所有时钟同步在约7 ms之内完成<sup>[41]</sup>。

借助时钟同步来处理分布式事务语义是一个非常有趣和活跃的研究领域<sup>[57,61,62]</sup>。但除了Google以外，目前主流数据库中还没有更多的实现。

---

注6： 还有一种分布式序列号生成器，例如Twitter的Snowflake，它用更为扩展的方式（例如将ID空间划分不同的范围，然后分配给不同的节点）来生成近似单调递增的唯一ID。但通常无法保证与因果关系一致的顺序。具体参阅第9章的“顺序保证”。



## 进程暂停

另一个分布式系统中危险使用时钟的例子：假设数据库每个分区只有一个主节点，只有主节点可以接受写入。那么其他节点该如何确信该主节点没有被宣告失效，可以安全地写入呢？

一种思路是主节点从其他节点获得一个租约，类似一个带有超时的锁<sup>[63]</sup>。某一个时间只有一个节点可以拿到租约，某节点获得租约之后，在租约到期之前，它就是这段时间内的主节点。为了维持主节点的身份，节点必须在到期之前定期去更新租约。如果节点发生了故障，则续约失败，这样另一个节点到期之后就可以接管。

典型处理流程如下所示：

```
while (true) {
    request = getIncomingRequest();

    // Ensure that the lease always has at least 10 seconds remaining
    if (lease.expiryTimeMillis - System.currentTimeMillis() < 10000) {
        lease = lease.renew();
    }

    if (lease.isValid()) {
        process(request);
    }
}
```

这代码有什么问题么？首先，它依赖于同步的时钟：租约到期时间由另一台机器所设置（例如，另一台机器的当前时间加上30秒得到到期时间），并和本地时钟进行比较。如果时钟之间有超过几秒的差异，这段代码会出现些奇怪的事情。

其次，即使我们改为仅使用本地单调时钟，还有另一个问题：代码假定时间检查点 `System.currentTimeMillis()` 与请求处理 `process(request)` 间隔很短，通常代码运行足够快，所以设置10秒的缓冲区来确保在请求处理过程中租约不会过期。

但是，如果程序执行中出现了某些意外的暂停呢？例如，假设线程在 `lease.isValid()` 消耗了整15秒。那么当开始处理请求时，租约已经过期，另一个节点已经接管了主节点。可惜我们无法有效通知线程暂停了这么长时间了，后续代码也不会注意到租约已经到期，除非运行到下一个循环迭代。不过，到那个时候它已经做了一些不安全的请求处理。

那么，一个线程可能会暂停这么长时间么？这是可能的，发生这种情况的原因有很多种：





- 许多编程语言（如Java虚拟机）都有垃圾收集器（GC），有时运行期间会暂停所有正在运行的线程。这些GC暂停甚至有时会持续数分钟<sup>[64]</sup>！即使像HotSpot JVM CMS所谓的“并发”垃圾收集器也不能完全与应用代码并行运行，需要时不时地停止活动的线程<sup>[65]</sup>。通过改变分配模式或调整GC某些参数可以减少一些暂停<sup>[66]</sup>，但是我们还是要防范最差情况以提供可靠的保证。
- 在虚拟化环境中，可能会暂停虚拟机（暂停所有执行进程并将内存状态保存到磁盘）然后继续（从内存中加载数据然后继续执行）。暂停可能发生进程运行的任一时刻，并且可能持续很长的时间。该功能通常用于实时迁移，即把虚拟机从一个主机迁移到另一个主机而不需要重启，这种情况下，暂停的长度主要取决于进程写入内存的速率<sup>[67]</sup>。
- 运行在终端用户设备（如笔记本电脑）时，执行也可能发生暂停，例如用户关闭了笔记本电脑或休眠。
- 当操作系统执行线程上下文切换时，或者虚拟机管理程序切换到另一个虚拟机时，正在运行的线程可能会在代码的任意位置被暂停。在虚拟机环境中，这种被其他虚拟机中断的CPU时间称为窃取时间。如果机器负载很高（即等待运行的线程很长），被暂停的线程可能需要一段时间之后才能再次运行。
- 如果应用程序执行同步磁盘操作，则线程可能暂停并等待磁盘I/O完成<sup>[68]</sup>。在许多语言中，即使代码并没有明确执行文件操作，也可能意外引入磁盘I/O。例如，Java类加载器在第一次使用类文件时会推迟加载，最终可能发生在执行时的任何时刻。I/O暂停和GC暂停甚至可能会同时发生，从而进一步恶化情况<sup>[69]</sup>。如果磁盘其实是个网络文件系统或网络块设备（如亚马逊的EBS），I/O还要受到网络延迟变化的影响<sup>[29]</sup>。
- 如果操作系统配置了基于磁盘的内存交换分区，内存访问可能触发缺页中断，进而需要从磁盘中加载内存页。I/O进行时（通常比较慢）线程为暂停。如果内存使用压力很大，还可能迫使更多的页面换出到磁盘。极端的情况下，操作系统可能会花费大量时间在页面换入换出上，而实际工作完成很少（所谓的颠簸）。为了避免此类问题，通常在服务器上禁用分页（宁愿杀死一些进程来释放内存而不是反复抖动）。
- 通过发送SIGSTOP信号来暂停UNIX进程，例如在shell中按下Ctrl-Z。这个信号会立即停止进程避免其拿到更多的CPU周期，直到接下来收到信号SIGCONT之后才从停止的地方继续运行。另外也不排除SIGSTOP信号是由运维人员不小心意外发送。



所有上述情况都可能随时抢占一个正在运行的线程，然后在之后的某个时间点再恢复线程的执行，而线程自身却对此一无所知。这个问题类似于在一台机器上运行多线程代码且保证线程安全。总之，你不能假定任何有关时间的东西，记住上下文切换和并行性可能随时可以发生。

在单台机器上编写多线程代码时，有不少工具可以帮助实现线程安全：互斥量，信号量，原子计数器，无锁数据结构，阻塞队列等。不幸的是，这些工具无法直接转为分布式系统，因为分布式系统通常不采用共享内存，而是在不可靠的网络上发送消息。

分布式系统中的一个节点必须假定，执行过程中的任何时刻都可能被暂停相当长一段时间，包括运行在某个函数中间。暂停期间，整个集群的其他部分都在照常运行，甚至一致将暂停的节点宣告为故障节点。最终，暂停的节点可能会回来继续运行，除非再次检查时钟，否则它对刚刚过去的暂停毫无意识。

## 响应时间保证

如上所述，在许多编程语言和操作系统中，线程和进程可能会暂停相当长的时间。如果仔细调教系统，可以做到避免很多这种暂停。

某些软件如果在指定时间内无法响应则会导致严重后果，这些运行环境包括：飞机，火箭，机器人，汽车和其他需要对输入传感器快速做出响应的组件等。对于这些系统，软件有一个必须做出响应的上限：如果无法满足，会导致系统级故障，这就是所谓的硬实时系统。



### 实时真的是实时吗？

在嵌入式系统中，实时通常意味着系统经过了精心设计和测试，以满足各种情况下执行时间约束。这与Web上模糊的“实时”术语的有着鲜明的对比，后者主要描述了一种持续的流式处理方式，但并没有强的时间约束。

例如，如果车载传感器检测到当前正在经历碰撞，肯定不希望系统由于不合适的GC暂停导致无法及时释放安全气囊。

提供实时保证需要来自软件栈的多个层面的支持：首先是一个实时操作系统（real-time operating system, RTOS），保证进程在给定的间隔内完成CPU时间片的调度分配；其次，库函数也必须考虑最坏的执行时间；然后，动态内存分配很可能要受限或者完全被禁止（如果存在实时垃圾收集器，确保GC不能处理太多任务）；最终还是需要大量、充分的测试和验证，以确保满足要求。



显然，这需要大量额外的工作，也严重限制了可使用的编程语言、库和工具的范围（因为大多数语言和工具库不提供实时保证）。由于这些原因，开发实时系统代价昂贵，通常只用于对安全至关重要的关键性嵌入式设备中。另外，“实时”与“高性能”不一样。实际上，实时系统往往吞吐量较低，它须优先考虑并响应高优先级的请求（参阅本章前面的“延迟与资源利用率”）。

对于大多数服务器端数据处理系统来说，实时性保证并不经济或者不合适。因此，现在这些运行在非实时环境下的系统就得承受如进程暂停、时钟不稳定等困扰。

## 调整垃圾回收的影响

无需昂贵的实时调度，还有一些其他措施可以减轻进程暂停所导致的负面影响。语言绑定的垃圾回收机制可以跟踪对象的分配情况以及剩余的空闲内存，因而可以在运行时灵活控制垃圾回收。

现在一个较新的想法是把GC暂停视为节点的一个计划内的临时离线，当节点启动垃圾回收时，通知其他节点来接管客户端的请求。此外，系统可以提前为前端应用发出预警，应用会等待当前请求完成，但停止向该节点发送新的请求，这样垃圾回收可以在无干扰的情况下更加高效运行。这个技巧以某种方式对客户端隐藏垃圾回收，降低负面影响<sup>[70,71]</sup>。目前一些对延迟敏感的系统（如金融交易系统<sup>[72]</sup>）已经采用了该方法。

该方法的一个变种是，只对短期对象（可以快速回收）执行垃圾回收，然后在其变成长期存活对象之前，采取定期重启的策略从而避免对长期存活对象执行全面回收<sup>[65,73]</sup>。每次选择一个节点重新启动，在重启之前，重新平衡节点之间的流量，思路与滚动升级类似（参阅第4章）。

这些措施虽然并不能完全避免垃圾回收导致的进程暂停，但可以有效地减少对应用层的影响。

## 知识，真相与谎言

本章到目前为止，已经探索了分布式系统与单节点程序的许多不同之处。例如，很少使用共享内存，通过不可靠网络传递消息且延迟不确定，可能遭受部分失效，不可靠的时钟以及进程暂停等。

如果还没有深谙分布式系统之道，那么下面些问题看起来有些难以理解。网络中的一个节点无法确信信息，只能通过网络收到（或没有收到）的消息来猜测。节点只能通过消息交换来获得其他节点当前的状态（存储了哪些数据，是否正常工作等）。如果



远程节点没有响应，由于没法区分网络本身的问题还是节点的问题，就无从知道节点究竟处于什么状态。

进一步深究上述问题，可能会涉及一些哲学相关话题：我们从系统获得的信息哪些是真实的、哪些是假的？如果感知和测量的手段都不可靠，那么获得的信息究竟有多大的可信度？软件系统是否应该遵循物理世界的那些通用法则（例如因果关系）呢？

还好，我们并不是非要搞清楚生命的意义才能回答上面的问题。在分布式系统中，我们可以明确列出对系统行为（系统模型）所做的若干假设，然后以满足这些假设条件为目标来构建实际运行的系统。在给定系统模型下，可以验证算法的正确性。这也意味着即使底层模型仅提供了少数几个保证，也可以在系统软件层面实现可靠的行为保证。

显然这不是一件容易的事情。在本章接下来的部分，我们将进一步探讨分布式系统中知识和真相，目的是帮忙我们审视可以做出哪些合理假设，以及通常可以提供哪些保证。第9章将介绍具体的例子，例如典型分布式系统和算法在特定条件下提供的哪些特定保证。

## 真相由多数决定

假定在一个发生了非对称故障的网络环境中，即某节点能够收到发送给它的消息，但是该节点发出的所有消息要么被丢弃，要么被推迟发送<sup>[19]</sup>。该节点即使本身运行良好，可以接收来自其他节点的请求，但其他节点却无法顺利收到响应。当消息超时之后，由于都收不到回复，其他节点就会一致声明上述节点发生失效。打个比方，就像这种情况：处于半连接的节点被强行摁倒在车里，即使它满心不愿意，即使在里面不停地哭喊“我没有问题”，可没人听到其呼救声，其他人只能忍痛将其塞进车里。

接下来是一个情况稍好的场景，半断开的节点可能会注意到其发送的消息没有被其他节点所确认，因此意识到网络一定发生了某种故障。尽管如此，节点还是会被其他节点错误地宣告为失效，改变不了该节点最终的命运。

第三种情况，该节点上垃圾回收运行了很长时间，所有线程包括那些事务处理任务都被GC抢占并暂停了足足一分钟，在此期间，没有处理任何请求，也没有发送任何响应。那么其他节点只能苦苦等待，不停重试，最后无奈宣布该节点已经失效，然后将其塞进车里。可节点最终还是完成了垃圾回收，原有的工作线程得以继续，好像什么都没有发生。此时，轮到其他节点感到惊讶：刚刚宣告有问题的节点突然从车里爬出来，活蹦乱跳，甚至和周围兴奋地聊天。显然，运行垃圾回收的节点根本没有意识它

中间休克了一分钟，从它的角度看，自上次正常的通信之后，一切正常，没有奇怪的停顿。

我们讲这些故事的寓意是，节点不能根据自己的信息来判断自身的状态。由于节点可能随时会失效，可能会暂停-假死，甚至最终无法恢复，因此，分布式系统不能完全依赖于单个节点。目前，许多分布式算法都依靠法定票数，即在节点之间进行投票（参阅第5章“读写quorum”）。任何决策都需要来自多个节点的最小投票数，从而减少对特定节点的依赖。

这其中包括关于宣告节点失效的决定。如果有法定数量的节点声明另一个节点失效，即使该节点仍感觉活得很自在，那它也必须接受失效的裁定，所有个体节点必须遵循法定投票的决议然后离线。

最常见的法定票数是取系统节点半数以上（也有其他类型的法定人数）。如果某些节点发生故障，quorum机制可以使系统继续工作（对于三个节点的系统，可以容忍一个节点失效；五个节点则可以容忍两个节点故障）。由于系统只可能存在一个多数，绝不会有多个多数在同时做出相互冲突的决定，因此系统的决议是可靠的。第9章介绍一致性算法时，我们将更详细地讨论如何使用quorum。

## 主节点与锁

有很多情况，我们需要在系统范围内只能有一个实例。例如：

- 只允许一个节点作为数据库分区的主节点，以防止出现脑裂（参阅第5章“处理节点失效”）。
- 只允许一个事务或客户端持有特定资源的锁，以防止同时写入从而导致数据破坏。
- 只允许一个用户来使用特定的用户名，从而确保用户名可以唯一标识用户。

在分布式系统实现时需要额外注意：即使某个节点自认为它是“唯一的那个”（例如分区的主节点，锁的持有者，成功拿走用户名的请求），但不一定获得了系统法定票数的同意！一个节点可能以前确实是主节点，但其他节点有可能在此期间已宣布其失效（例如，出现了网络中断或GC暂停），节点已被降级而系统选出了另一个主节点。

当多数节点声明节点已失效，而该节点还继续充当“唯一的那个”，如果系统设计不周就会导致负面后果。该节点会按照自认为正确的信息向其他节点发送消息，其他节点如果还选择相信它，那么系统就会出现错误的行为。

例如图8-4展示了由于不正确的加锁而导致数据破坏的例子。该bug并非只是理论存在，HBase曾遭遇过该问题<sup>[74,75]</sup>。其设计目标是确存储系统的文件一次只能由一个客户端访问，如果多个客户端试图同时写入该文件，文件就会被破坏。因此，在访问文件之前客户端需要从锁服务获取访问租约。

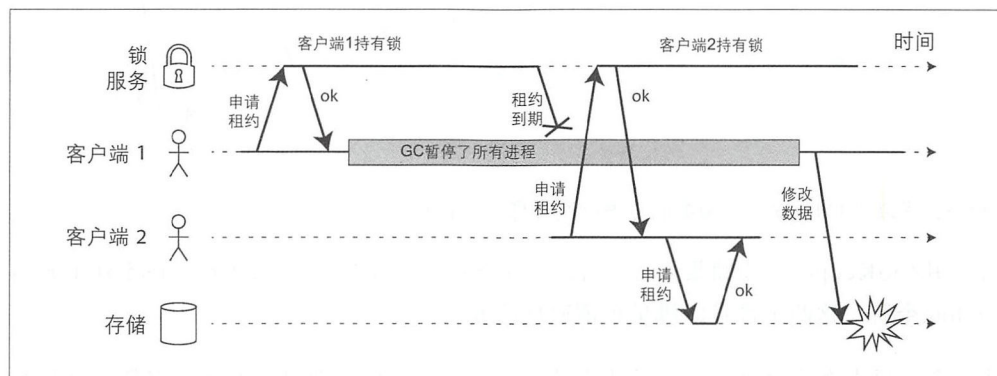


图8-4：分布式锁的不正确实现：客户端1的锁租约其实已经过期，但它自认为有效，最终导致文件破坏

这个问题属于前面“进程暂停”中的一种情况：持有租约的客户端被暂停太久直到租约到期。然后另一个客户端已经获得了文件的锁租约，并开始写文件。接下来，当暂停的客户端重新回来时，它仍然（错误地）认为合法持有锁并尝试写文件。结果导致客户2的文件写入被破坏。

## Fencing令牌

当使用锁和租约机制来保护资源的并发访问时（见图8-4），必须确保过期的“唯一的那个”节点不能影响其他正常部分。要实现这一目标，可以采用一种相当简单的技术 *fencing*（栅栏，隔离之意），如图8-5所示。

我们假设每次锁服务在授予锁或租约时，还会同时返回一个 *fencing* 令牌，该令牌（数字）每授予一次就会递增（例如，由锁服务增加）。然后，要求客户端每次向存储服务发送写请求时，都必须包含所持有的 *fencing* 令牌。

图8-5中，客户端1获得锁租约的同时得到了令牌号33，但随后陷入了一个长时间的暂停直到租约到期。这时客户端2已经获得了锁租约和令牌号34，然后发送写请求（以及令牌号34）到存储服务。接下来客户端1恢复过来，并以令牌号33来尝试写入，存储服务由于记录了最近已经完成了更高令牌号（34），因此拒绝令牌号33的写请求。



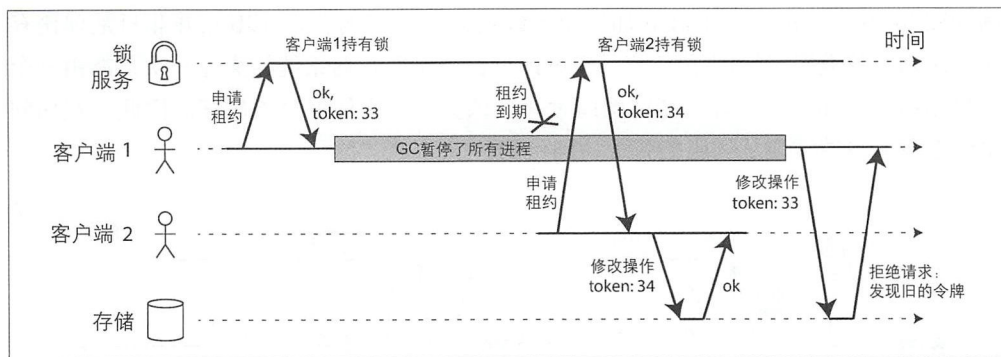


图8-5：递增fencing令牌，拒绝旧令牌的写操作以确保存储安全

当使用ZooKeeper作为锁服务时，可以用事务标识`zxid`或节点版本`cversion`来充当fencing令牌，这两个都可以满足单调递增的要求<sup>[74]</sup>。

请注意，只靠客户端自己检查锁状态是不够的，这种机制要求资源本身必须主动检查所持令牌信息，如果发现已经处理过更高令牌请求，要拒绝持有低令牌的所有写请求。如果资源不支持额外的令牌检查，可以采取一些临时技巧来绕过去（例如，对于访问文件存储服务的情况，可以将令牌信息内嵌在文件名中）。总之，为了避免在锁保护之外发生请求处理，需要进行额外的检查机制。

在服务器端检查令牌可能看起来有些复杂，但其实是推荐的正确做法：系统服务不能假定所有的客户端都表现符合预期，事实上客户端通常由权限级别相对较低的人来操作运行<sup>[76]</sup>，因此存在一定的误用、滥用风险，从安全角度讲，服务端必须防范这种来自客户端的滥用。

## 拜占庭故障

fencing令牌可以检测并阻止那些无意的误操作（例如节点并没有发现其租约已经过期）。但是，如果节点故意试图破坏系统，在发送消息时可以简单地伪造令牌即可。

本书总是假设节点虽然不可靠但一定是诚实的：它们尽管运行很慢或者由于故障而无法响应，或者状态可能已经过期（例如由于GC暂停或网络延迟），但一旦做出了响应，则一定是完全基于其所知的全部信息和事先协议约定好的行为准则，响应代表了其所知的“真相”。

如果节点存在“撒谎”的情况（即故意发送错误的或破坏性的响应），那么分布式系统处理的难度就上了一个台阶。例如，节点明明没有收到某条消息，但却对外声称收

到了。这种行为称为拜占庭故障，在这样不信任的环境中需要达成共识的问题也被称为拜占庭将军问题<sup>[77]</sup>。

## 拜占庭将军问题

拜占庭将军问题是所谓“两将军问题”的更抽象表示<sup>[78]</sup>，后者假定有两名将军需要就战斗计划达成一致。由于他们在两个不同的地点建立了营地，中间只能通过信使进行沟通，而信使在传递消息时可能会出现延迟或丢失（就像网络中的信息包一样）。我们将在第9章讨论共识问题。

而在拜占庭版本中，有 $n$ 位将军需要达成共识，并且其中存在一些叛徒试图阻挠达成共识。大多数的将军都是忠诚的，发出了真实的信息，但是叛徒则试图通过发送虚假或不真实的信息来欺骗和混淆他人（同时努力隐藏自己）。而且大家事先并不知道叛徒是谁。

拜占庭是一个古希腊城市，即随后的君士坦丁堡，目前位于土耳其的伊斯坦布尔。没有任何历史证据表明，拜占庭将军比其他地方的将军更诡计多端或者善于密谋。这个名字其实隐喻着拜占庭式的过分复杂，官僚，含混不清的意思，早在计算机之前已在政治中广泛引用<sup>[79]</sup>。拜占庭问题的提出者兰波特当初只是想选一个不会冒犯任何人/国籍的名字，据传他曾被明确通牒，不准使用阿尔巴尼亚将军问题来命名<sup>[80]</sup>。

如果某个系统中即使发生部分节点故障，甚至不遵从协议，或者恶意攻击、干扰网络，但仍可继续正常运行，那么我们称之为拜占庭式容错系统。这些担忧在某些特定场景是合理的。例如：

- 在航空航天领域，计算机内存或CPU寄存器中的数据可能会被辐射而发生故障，导致以不可预知的方式响应其他节点。这种情况下如果将系统下线，代价将异常昂贵（例如，可能出现飞机撞毁，杀死船员，或致使火箭与国际空间站相撞等），飞行控制系统必须做到容忍拜占庭故障<sup>[81,82]</sup>。
- 在有多个参与者的系统中，某些参与者可能会作弊或者欺骗他人。这时节点不能完全相信另一个节点所发送的消息，它可能就是恶意的。例如，像比特币和其他区块链一样的点对点网络就是让互不信任的当事方就某项交易达成一致，且不依赖于集中的机制<sup>[83]</sup>。

然而，在本书所讨论的这些系统中，我们可以安全地假定没有拜占庭式的故障。在数

据中心里，所有的节点都是由一个组织来集中控制（可信任），辐射水平也足够低因而内存损坏可以忽略。解决拜占庭容错的系统协议异常复杂<sup>[84]</sup>，而容错的嵌入式系统还依赖于硬件层面的支持<sup>[81]</sup>。因而在绝大多数服务器端数据系统中，部署拜占庭容错解决方案基本不太可行。

Web应用程序确实可能会接收来自任意客户端（如Web浏览器）的请求，其中可能带有恶意行为。因此需要输入验证、安全监测和输出转义等步骤，例如防止SQL注入和跨站恶意脚本。但我们通常并不使用拜占庭容错协议，而只是全权让服务器决定什么是可接受的客户端行为，什么是不允许的。只有在没有这种中央决策机制的点对点网络中，拜占庭容错才更为必要。

另外，软件中的bug可以被认为是拜占庭式故障，但如果将相同的软件部署到所有节点上，那么即使拜占庭式的容错算法也无法解决问题。大多数拜占庭容错算法要求系统超过三分之二的节点即绝大多数要功能正常（如果有四个节点，则最多允许一台发生故障）。要采用这类算法对付bug，必须有四个不同的软件实现，然后寄希望该bug只出现在四个实现中的一个。

试想如果有这样一个协议或者算法能同时保护我们免受漏洞，安全防范降级以及恶意攻击等，那将是多么美好。不幸的是，这是不现实的。通常如果攻击者可以入侵一个节点，则很可能会攻陷几乎所有节点（由于运行相同的软件）。因此，传统的安全措施如认证、访问控制、加密、防火墙等，仍是防范攻击的主要保护机制。

## 弱的谎言形式

尽管我们假设节点通常是诚实的，但依然推荐增加必要的机制来防范一些不那么恶意的“谎言”。例如由于硬件问题造成的无效消息、软件bug和配置错误。这种保护机制显然并不是完整的拜占庭式容错，无法防范敌手的攻击，但它们更为简单实用，可以帮助提高软件系统的可靠性和健壮性。例如：

- 由于硬件问题或操作系统、驱动程序、路由器等方面的错误，导致网络数据包有时出现损坏。通常，可以借助TCP/UDP中内置的数据包校验和来发现这类问题，但有时他们会逃避检测<sup>[85-87]</sup>。此时，一个简单的防范措施是在应用层添加校验和。
- 对公众开放的应用必须仔细检查用户的所有输入，例如输入值是否在合理的范围内，并限制字符串的大小，防止分配超大内存导致拒绝服务攻击。位于防火墙后面的内部服务可能不用太严格的检查，但依然推荐执行基本的安全检查（例如，接口协议解析<sup>[85]</sup>）。



- NTP客户端最好配置多个时间服务器。同步时间时，连接到多个时间服务器，收到回应之后要评估时间偏差，使得多数服务器就一定的时间范围达成一致。只要大多数服务器都正常，某个错误的服务器就可以被检测出来，从而排除在同步结果之外<sup>[37]</sup>。总之，使用多台NTP服务器可以比仅使用一台服务器更为鲁棒。

## 理论系统模型与现实

目前分布式系统方面已有许多不错的具体算法，如第9章中要介绍的共识算法。这些算法需要容忍本章所讨论的各种故障。

算法的实现不能过分依赖特定的硬件和软件配置。这就要求我们需要对预期的系统错误进行形式化描述。我们通过定义一些系统模型来形式化描述算法的前提条件。

关于计时方面，有三种常见的系统模型：

### 同步模型

同步模型假定有上界的网络延迟，有上界的进程暂停和有上界的时钟误差。注意，这并不意味着完全同步的时钟或者网络延迟为零。它只意味着你清楚地了解网络延迟、暂停和时钟漂移不会超过某个固定的上限<sup>[88]</sup>。大多数实际系统的实际模型并非同步模型，因为（如本章所讨论的）无限延迟和暂停确实可能发生。

### 部分同步模型

部分同步意味着系统在大多数情况下像一个同步系统一样运行，但有时候会超出网络延迟，进程暂停和时钟漂移的预期上界<sup>[88]</sup>。这是一个比较现实的模型：大多数情况下，网络和进程比较稳定（否则几乎不可能提供持续的服务），但是我们必须考虑到任何关于时机的假设都有偶尔违背的情况，而一旦发生，网络延迟，暂停和时钟偏差可能会变得非常大。

### 异步模型

在这个模型中，一个算法不会对时机做任何的假设，甚至里面根本没有时钟（也就没有超时机制）。某些算法可以支持纯异步模型，但并不常见。

除了时机之外，我们还需要考虑节点失效。有以下三种最常见的节点失效系统模型：

### 崩溃-中止模型

在崩溃-中止模型中，算法假设一个节点只能以一种方式发生故障，即遭遇系统崩溃。这意味着节点可能在任何时候突然停止响应，且该节点以后永远消失，无法恢复。

## 崩溃-恢复模型

节点可能会在任何时候发生崩溃，且可能会在一段（未知的）时间之后得到恢复并再次响应。在崩溃-恢复模型中，节点上持久性存储（即非易失性存储）的数据会在崩溃之后得以保存，而内存中状态可能会丢失。

## 拜占庭（任意）失效模型

如上节所述，节点可能发生任何事情，包括试图作弊和欺骗其他节点。

对于真实系统的建模，最普遍的组合是崩溃-恢复模型结合部分同步模型。那么接下来上层的分布式算法该如何应对这样的模型呢？

## 算法的正确性

为了定义算法的正确性，我们可以描述它的属性信息。例如，排序算法的输出具有以下特性：对于输出列表中的任何两个不同的元素，左边的元素小于右边的元素。这就是一个对列表进行排序的正确性描述。

类似的思路，我们可以通过描述目标分布式算法的相关属性来定义其正确性。例如，对于锁服务的fencing 令牌生成算法（参阅本章前面的“fencing令牌”），要求算法具有以下属性：

### 唯一性

两个令牌请求不能获得相同的值。

### 单调递增

如果请求 $x$ 返回了令牌 $t_x$ ，请求 $y$ 返回了令牌 $t_y$ ，且 $x$ 在 $y$ 开始之前先完成，那么 $t_x < t_y$ 。

### 可用性

请求令牌的节点如果不发生崩溃则最终一定会收到响应。

如果针对某个系统模型的算法在各种情况下都能满足定义好的属性要求，那么我们称这个算法是正确的。这有何意义呢？如果换一个角度来看一种极端情况，所有节点全部崩溃，或者所有的网络延迟突然变得无限长，那么所有的算法都不可能完成其预期功能。

## 安全与活性

为进一步加深理解，有必要区分两种不同的属性：安全性和活性。在上面的例子中，唯一性和单调递增属于安全属性，而可用性则属于活性。



这两种性质有何区别呢？一种理解思路是，活性的定义中通常会包括暗示“最终”一词（是的，你猜对了，最终一致性也是一种活性<sup>[89]</sup>）。

安全性通常可以理解为“没有发生意外”，而活性则类似“预期的事情最终一定会发生”。这个非正式定义中，有很多主观因素，所以不用过度解读。安全性和活性其实有准确和数学化的定义描述，请参阅文献<sup>[90]</sup>：

- 如果违反了安全属性，我们可以明确指向发生的特定的时间点（例如，唯一性如果被违反，我们可以定位到具体哪个操作产生了重复令牌）。且一旦违反安全属性，违规行为无法撤销，破坏已实际发生。
- 活性则反过来：可能无法明确某个具体的时间点（例如一个节点发送了一个请求，但还没有收到响应），但总是希望在未来某个时间点可以满足要求（即收到回复）。

区分安全性和活性的一个好处是可以帮助简化处理一些具有挑战性的系统模型。通常对于分布式算法，要求在所有可能的系统模型下，都必须符合安全属性<sup>[88]</sup>。也就是说，即使所有节点发生崩溃，或者整个网络中断，算法确保不会返回错误的结果。

而对于活性，则存在一些必要条件。例如，我们可以说，只有在多数节点没有崩溃，以及网络最终可以恢复的前提下，我们才能保证最终可以收到响应。部分同步模型的定义即要求任何网络中断只会持续一段有限的时间，然后得到了修复，系统最终返回到同步的一致状态。

## 将系统模型映射到现实世界

安全性、活性以及所建立的系统模型对于评测分布式算法的正确性意义重大。然而，很明显系统模型只是对现实情况的简化抽象，实践中具体实施算法时，各种因素混杂在一起会提出更严峻的挑战。

例如，在崩溃-恢复模型中，算法通常假设保存在持久性介质的数据可以安然无恙。但是，如果硬盘上的数据发生损坏，或者由于硬件错误，配置错误<sup>[91]</sup>等导致数据被清除，会发生什么后果呢？即使硬盘本身正确连接到服务器，但服务器存在固件错误，导致重启时无法正确识别硬盘，又会发生什么情况？

Quorum算法（参见第5章“读写Quorum”）要求节点必须记录之前对外所宣告的数据。如果节点发生意外而丢弃存储的数据，会打破法定条件并破坏算法的正确性。或许此时我们需要一个新的系统模型，它假定通常情况下数据存储非常可靠，但还是有丢失的可能。但那个模型就变得更难以推理了。





算法的理论描述可以简单地宣称某些事情绝不会发生，例如在非拜占庭系统中，我们就会明确假定一些可能和不可能发生的错误。不过，真正去实现时最好还是有一些必要的代码来简单处理一些几乎不可能发生的事情，即使只是去输出一些提醒信息（如 `printf`）和程序退出错误代码（如 `exit-666`），以方便操作人员来清理最后的烂摊子<sup>[93]</sup>。这些错误处理也很好地体现了计算机科学和软件工程之间的差异。

这绝对不是说抽象的系统模型没有价值，恰恰相反，它把实际系统中的复杂性提炼成一个更容易理解、更具可控性的抽象错误集合，可以有效帮助我们理解问题之本质，然后设计系统性方法来最终解决问题。如果我们可以确定在给定的系统模型中，算法总能满足属性要求，那么我们可以证明算法就是正确的。

证明算法正确并不意味着真实系统上的某个具体实现一定是正确的。毫无疑问，这依然是极其重要的第一步。理论分析并不能覆盖现实系统隐藏的每一个问题细节，例如某些边界条件一旦触发，会有负面影响，但这些边界条件归根结底一定是违背了模型的某个前提假设（例如，关于时间点）。我们可以说，理论性分析与实证性检验对最终的成功同等重要。

## 小结

本章讨论了分布式系统中可能发生各种典型问题，包括：

- 当通过网络发送数据包时，数据包可能会丢失或者延迟；同样，回复也可能会丢失或延迟。所以如果没有收到回复，并不能确定消息是否发送成功。
- 节点的时钟可能会与其他节点存在明显的不同步（尽管尽最大努力设置了NTP服务器），时钟还可能会突然向前跳跃或者倒退，依靠精确的时钟存在一些风险，没有特别简单的办法来精确测量时钟的偏差范围。
- 进程可能在执行过程中的任意时候遭遇长度未知的暂停（一个重要的原因是垃圾回收），结果它被其他节点宣告为失效，尽管后来又恢复执行，却对中间的暂停毫无所知。

部分失效可能是分布式系统的关键特征。只要软件试图跨节点做任何事情，就有可能出现失败，或者随机变慢，或者根本无应答（最终超时）。对于分布式环境，我们的目标是建立容忍部分失效的软件系统，这样即使某些部件发生失效，系统整体还可以继续运行。

为了容忍错误，第一步是检测错误，但即使这样也很有挑战。多数系统没有检测节点是否发生故障的准确机制，因此分布式算法更多依靠超时来确定远程节点是否仍然可



用。但是，超时无法区分网络和节点故障，且可变的网络延迟有时会导致节点被误认为发生崩溃。此外，节点可能处于一种降级状态：例如，由于驱动程序错误，千兆网络接口可能突然降到1kb/s的吞吐量<sup>[94]</sup>。这样一个处于“残废”的节点比彻底挂掉的故障节点更难处理。

检测到错误之后，让系统容忍失效也不容易。在典型的分布式环境下，没有全局变量，没有共享内存，没有约定的尝试或其他跨节点的共享状态。节点甚至不太清楚现在的准确时间，更不用说其他更高级的了。信息从一个节点流动到另一个节点只能是通过不可靠的网络来发送。单个节点无法安全的做出任何决策，而是需要多个节点之间的共识协议，并争取达到法定票数。

如果习惯于编写单节点理想化环境运行的软件（即同一个操作总是确定性地返回相同的结果），当转向分布式系统时，种种看似凌乱的现实可能着实让人震惊。相反，如果在单节点上即可解决问题，那么对于一个分布式系统工程师通常会被认为该问题微不足道（现在单节点确实可以完成很多任务<sup>[95]</sup>）。如果可以避免打开潘多拉之盒，那么把工作都放在一台机器也值得一试。

但正如在第二部分开头所讨论的那样，可扩展性并不是使用分布式系统的唯一原因。容错与低延迟（将数据放置在距离用户较近的地方）也是同样重要的目标，而后两者无法靠单节点来实现。

本章，我们也探讨了网络、时钟和进程的不可靠性是否是不可避免的自然规律。我们对此给出的结论是否定的：的确有可能在网络中提供硬实时的延迟保证或者具有上确界的延迟，但代价昂贵，且硬件资源利用率很低。除了安全关键场景，目前绝大多数都选择了低成本（和不可靠）。

我们还谈到了高性能计算，它们多采用更加可靠的组件，发生故障时完全停止系统，之后重新启动。相比之下，分布式系统会长时间不间断运行，以避免影响服务级别；故障处理和系统维护多以节点为单位进行处理，或者理论上如此（实际上，如果错误的配置不小心被应用到集群的所有节点，仍然会导致整个集群系统瘫痪）。

这样看起来本章揭露的全是问题，前景黯淡。那么在下一章，我们将讨论解决方案，重点是针对这些问题而设计相关的分布式算法。

## 参考文献

[1] Mark Cavage: “There’s Just No Getting Around It: You’re Building a Distributed System,” *ACM Queue*, volume 11, number 4, pages 80-89, April 2013. doi:



10.1145/2466486.2482856.

[2] Jay Kreps: “Getting Real About Distributed System Reliability,” *blog.empathybox.com*, March 19, 2012.

[3] Sydney Padua: *The Thrilling Adventures of Lovelace and Babbage: The (Mostly) True Story of the First Computer*. Particular Books, April 2015. ISBN: 978-0-141-98151-2.

[4] Coda Hale: “You Can’t Sacrifice Partition Tolerance,” *codahale.com*, October 7, 2010.

[5] Jeff Hodges: “Notes on Distributed Systems for Young Bloods,” *somethingsimilar.com*, January 14, 2013.

[6] Antonio Regalado: “Who Coined ‘Cloud Computing’?,” *technologyreview.com*, October 31, 2011.

[7] Luiz André Barroso, Jimmy Clidaras, and Urs Hölzle: “The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines, Second Edition,” *Synthesis Lectures on Computer Architecture*, volume 8, number 3, Morgan & Claypool Publishers, July 2013. doi:10.2200/S00516ED2V01Y201306CAC024, ISBN: 978-1-627-05010-4.

[8] David Fiala, Frank Mueller, Christian Engelmann, et al.: “Detection and Correction of Silent Data Corruption for Large-Scale High-Performance Computing,” at *International Conference for High Performance Computing, Networking, Storage and Analysis (SC12)*, November 2012.

[9] Arjun Singh, Joon Ong, Amit Agarwal, et al.: “Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google’s Datacenter Network,” at *Annual Conference of the ACM Special Interest Group on Data Communication (SIGCOMM)*, August 2015. doi:10.1145/2785956.2787508.

[10] Glenn K. Lockwood: “Hadoop’s Uncomfortable Fit in HPC,” *glennklockwood.blogspot.co.uk*, May 16, 2014.

[11] John von Neumann: “Probabilistic Logics and the Synthesis of Reliable Organisms from Unreliable Components,” in *Automata Studies (AM-34)*, edited by Claude E. Shannon and John McCarthy, Princeton University Press, 1956. ISBN: 978-0-691-07916-5.





- [12] Richard W. Hamming: *The Art of Doing Science and Engineering*. Taylor & Francis, 1997. ISBN: 978-9-056-99500-3.
- [13] Claude E. Shannon: “A Mathematical Theory of Communication,” *The Bell System Technical Journal*, volume 27, number 3, pages 379-423 and 623-656, July 1948.
- [14] Peter Bailis and Kyle Kingsbury: “The Network Is Reliable,” *ACM Queue*, volume 12, number 7, pages 48-55, July 2014. doi:10.1145/2639988.2639988.
- [15] Joshua B. Leners, Trinabh Gupta, Marcos K. Aguilera, and Michael Walfish: “Taming Uncertainty in Distributed Systems with Help from the Network,” at *10th European Conference on Computer Systems (EuroSys)*, April 2015. doi:10.1145/2741948.2741976.
- [16] Phillipa Gill, Navendu Jain, and Nachiappan Nagappan: “Understanding Network Failures in Data Centers: Measurement, Analysis, and Implications,” at *ACM SIGCOMM Conference*, August 2011. doi:10.1145/2018436.2018477.
- [17] Mark Imbriaco: “Downtime Last Saturday,” *github.com*, December 26, 2012.
- [18] Will Oremus: “The Global Internet Is Being Attacked by Sharks, Google Confirms,” *slate.com*, August 15, 2014.
- [19] Marc A. Donges: “Re: bnx2 cards Intermittantly Going Offline,” Message to Linux *netdev* mailing list, spinics.net, September 13, 2012.
- [20] Kyle Kingsbury: “Call Me Maybe: Elasticsearch,” *aphyr.com*, June 15, 2014.
- [21] Salvatore Sanfilippo: “A Few Arguments About Redis Sentinel Properties and Fail Scenarios,” *antirez.com*, October 21, 2014.
- [22] Bert Hubert: “The Ultimate SO\_LINGER Page, or: Why Is My TCP Not Reliable,” *blog.netherlabs.nl*, January 18, 2009.
- [23] Nicolas Liochon: “CAP: If All You Have Is a Timeout, Everything Looks Like a Partition,” *blog.thislongrun.com*, May 25, 2015.
- [24] Jerome H. Saltzer, David P. Reed, and David D. Clark: “End-To-End Arguments in System Design,” *ACM Transactions on Computer Systems*, volume 2, number 4, pages 277-288, November 1984. doi:10.1145/357401.357402.



- [25] Matthew P. Grosvenor, Malte Schwarzkopf, Ionel Gog, et al.: “Queues Don’ t Matter When You Can JUMP Them!,” at *12th USENIX Symposium on Networked Systems Design and Implementation* (NSDI), May 2015.
- [26] Guohui Wang and T. S. Eugene Ng: “The Impact of Virtualization on Network Performance of Amazon EC2 Data Center,” at *29th IEEE International Conference on Computer Communications* (INFOCOM), March 2010. doi:10.1109/INFCOM.2010.5461931.
- [27] Van Jacobson: “Congestion Avoidance and Control,” at *ACM Symposium on Communications Architectures and Protocols* (SIGCOMM), August 1988. doi: 10.1145/52324.52356.
- [28] Brandon Philips: “etcd: Distributed Locking and Service Discovery,” at *Strange Loop*, September 2014.
- [29] Steve Newman: “A Systematic Look at EC2 I/O,” *blog.scalyr.com*, October 16, 2012.
- [30] Naohiro Hayashibara, Xavier Défago, Rami Yared, and Takuya Katayama: “The  $\phi$  Accrual Failure Detector,” Japan Advanced Institute of Science and Technology, School of Information Science, Technical Report IS-RR-2004-010, May 2004.
- [31] Jeffrey Wang: “Phi Accrual Failure Detector,” *ternarysearch.blogspot.co.uk*, August 11, 2013.
- [32] Srinivasan Keshav: *An Engineering Approach to Computer Networking: ATM Networks, the Internet, and the Telephone Network*. Addison-Wesley Professional, May 1997. ISBN: 978-0-201-63442-6.
- [33] Cisco, “Integrated Services Digital Network,” *docwiki.cisco.com*.
- [34] Othmar Kyas: *ATM Networks*. International Thomson Publishing, 1995. ISBN: 978-1-850-32128-6.
- [35] “InfiniBand FAQ,” Mellanox Technologies, December 22, 2014.
- [36] Jose Renato Santos, Yoshio Turner, and G. (John) Janakiraman: “End-to-End Congestion Control for InfiniBand,” at *22nd Annual Joint Conference of the IEEE Computer and Communications Societies* (INFOCOM), April 2003. Also published



by HP Laboratories Palo Alto, Tech Report HPL-2002-359. doi:10.1109/INFCOM.2003.1208949.

[37] Ulrich Windl, David Dalton, Marc Martinec, and Dale R. Worley: “The NTP FAQ and HOWTO,” *ntp.org*, November 2006.

[38] John Graham-Cumming: “How and why the leap second affected Cloudflare DNS,” *blog.cloudflare.com*, January 1, 2017.

[39] David Holmes: “Inside the Hotspot VM: Clocks, Timers and Scheduling Events - Part I - Windows,” *blogs.oracle.com*, October 2, 2006.

[40] Steve Loughran: “Time on Multi-Core, Multi-Socket Servers,” *stevelloughran.blogspot.co.uk*, September 17, 2015.

[41] James C. Corbett, Jeffrey Dean, Michael Epstein, et al.: “Spanner: Google’s Globally-Distributed Database,” at *10th USENIX Symposium on Operating System Design and Implementation (OSDI)*, October 2012.

[42] M. Caporali and R. Ambrosini: “How Closely Can a Personal Computer Clock Track the UTC Timescale Via the Internet?,” *European Journal of Physics*, volume 23, number 4, pages L17-L21, June 2012. doi:10.1088/0143-0807/23/4/103.

[43] Nelson Minar: “A Survey of the NTP Network,” *alumni.media.mit.edu*, December 1999.

[44] Viliam Holub: “Synchronizing Clocks in a Cassandra Cluster Pt. 1-The Problem,” *blog.logentries.com*, March 14, 2014.

[45] Poul-Henning Kamp: “The One-Second War (What Time Will You Die?),” *ACM Queue*, volume 9, number 4, pages 44-48, April 2011. doi: 10.1145/1966989.1967009.

[46] Nelson Minar: “Leap Second Crashes Half the Internet,” *somebits.com*, July 3, 2012.

[47] Christopher Pascoe: “Time, Technology and Leaping Seconds,” *googleblog.blogspot.co.uk*, September 15, 2011.

[48] Mingxue Zhao and Jeff Barr: “Look Before You Leap - The Coming Leap Second and AWS,” *aws.amazon.com*, May 18, 2015.



- [49] Darryl Veitch and Kanthaiyah Vijayalayan: “Network Timing and the 2015 Leap Second,” at *17th International Conference on Passive and Active Measurement (PAM)*, April 2016. doi:10.1007/978-3-319-30505-9\_29.
- [50] “Timekeeping in VMware Virtual Machines,” Information Guide, VMware, Inc., December 2011.
- [51] “MiFID II / MiFIR: Regulatory Technical and Implementing Standards-Annex I (Draft),” European Securities and Markets Authority, Report ESMA/2015/1464, September 2015.
- [52] Luke Bigum: “Solving MiFID II Clock Synchronisation With Minimum Spend (Part 1),” *lmax.com*, November 27, 2015.
- [53] Kyle Kingsbury: “Call Me Maybe: Cassandra,” *aphyr.com*, September 24, 2013.
- [54] John Daily: “Clocks Are Bad, or, Welcome to the Wonderful World of Distributed Systems,” *basho.com*, November 12, 2013.
- [55] Kyle Kingsbury: “The Trouble with Timestamps,” *aphyr.com*, October 12, 2013.
- [56] Leslie Lamport: “Time, Clocks, and the Ordering of Events in a Distributed System,” *Communications of the ACM*, volume 21, number 7, pages 558-565, July 1978. doi:10.1145/359545.359563.
- [57] Sandeep Kulkarni, Murat Demirbas, Deepak Madeppa, et al.: “Logical Physical Clocks and Consistent Snapshots in Globally Distributed Databases,” State University of New York at Buffalo, Computer Science and Engineering Technical Report 2014-04, May 2014.
- [58] Justin Sheehy: “There Is No Now: Problems With Simultaneity in Distributed Systems,” *ACM Queue*, volume 13, number 3, pages 36-41, March 2015. doi: 10.1145/2733108.
- [59] Murat Demirbas: “Spanner: Google’s Globally-Distributed Database,” muratbuffalo. *blogspot.co.uk*, July 4, 2013.
- [60] Dahlia Malkhi and Jean-Philippe Martin: “Spanner’s Concurrency Control,” *ACM SIGACT News*, volume 44, number 3, pages 73-77, September 2013. doi: 10.1145/2527748.2527767.

- [61] Manuel Bravo, Nuno Diegues, Jingna Zeng, et al.: “On the Use of Clocks to Enforce Consistency in the Cloud,” *IEEE Data Engineering Bulletin*, volume 38, number 1, pages 18-31, March 2015.
- [62] Spencer Kimball: “Living Without Atomic Clocks,” *cockroachlabs.com*, February 17, 2016.
- [63] Cary G. Gray and David R. Cheriton: “Leases: An Efficient Fault-Tolerant Mechanism for Distributed File Cache Consistency,” at *12th ACM Symposium on Operating Systems Principles (SOSP)*, December 1989. doi:10.1145/74850.74870.
- [64] Todd Lipcon: “Avoiding Full GCs in Apache HBase with MemStore-Local Allocation Buffers: Part 1,” *blog.cloudera.com*, February 24, 2011.
- [65] Martin Thompson: “Java Garbage Collection Distilled,” *mechanicalsympathy.blogspot.co.uk*, July 16, 2013.
- [66] Alexey Ragozin: “How to Tame Java GC Pauses? Surviving 16GiB Heap and Greater,” *java.dzone.com*, June 28, 2011.
- [67] Christopher Clark, Keir Fraser, Steven Hand, et al.: “Live Migration of Virtual Machines,” at *2nd USENIX Symposium on Symposium on Networked Systems Design & Implementation (NSDI)*, May 2005.
- [68] Mike Shaver: “fsyncers and Curveballs,” *shaver.off.net*, May 25, 2008.
- [69] Zhenyun Zhuang and Cuong Tran: “Eliminating Large JVM GC Pauses Caused by Background IO Traffic,” *engineering.linkedin.com*, February 10, 2016.
- [70] David Terei and Amit Levy: “Blade: A Data Center Garbage Collector,” arXiv: 1504.02578, April 13, 2015.
- [71] Martin Maas, Tim Harris, Krste Asanović, and John Kubiatowicz: “Trash Day: Coordinating Garbage Collection in Distributed Systems,” at *15th USENIX Workshop on Hot Topics in Operating Systems (HotOS)*, May 2015.
- [72] “Predictable Low Latency,” Cinnober Financial Technology AB, *cinnober.com*, November 24, 2013.
- [73] Martin Fowler: “The LMAX Architecture,” *martinfowler.com*, July 12, 2011.

[74] Flavio P. Junqueira and Benjamin Reed: *ZooKeeper: Distributed Process Coordination*. O'Reilly Media, 2013. ISBN: 978-1-449-36130-3.

[75] Enis Söztutar: “HBase and HDFS: Understanding Filesystem Usage in HBase,” at *HBaseCon*, June 2013.

[76] Caitie McCaffrey: “Clients Are Jerks: AKA How Halo 4 DoSed the Services at Launch & How We Survived,” *caitiem.com*, June 23, 2015.

[77] Leslie Lamport, Robert Shostak, and Marshall Pease: “The Byzantine Generals Problem,” *ACM Transactions on Programming Languages and Systems (TOPLAS)*, volume 4, number 3, pages 382-401, July 1982. doi:10.1145/357172.357176.

[78] Jim N. Gray: “Notes on Data Base Operating Systems,” in *Operating Systems: An Advanced Course*, Lecture Notes in Computer Science, volume 60, edited by R. Bayer, R. M. Graham, and G. Seegmüller, pages 393-481, Springer-Verlag, 1978. ISBN: 978-3-540-08755-7.

[79] Brian Palmer: “How Complicated Was the Byzantine Empire?,” *slate.com*, October 20, 2011.

[80] Leslie Lamport: “My Writings,” *research.microsoft.com*, December 16, 2014. This page can be found by searching the web for the 23-character string obtained by removing the hyphens from the string `allla-mport-spubso-ntheweb`.

[81] John Rushby: “Bus Architectures for Safety-Critical Embedded Systems,” at *1st International Workshop on Embedded Software (EMSOFT)*, October 2001.

[82] Jake Edge: “ELC: SpaceX Lessons Learned,” *lwn.net*, March 6, 2013.

[83] Andrew Miller and Joseph J. LaViola, Jr.: “Anonymous Byzantine Consensus from Moderately-Hard Puzzles: A Model for Bitcoin,” University of Central Florida, Technical Report CS-TR-14-01, April 2014.

[84] James Mickens: “The Saddest Moment,” *USENIX ;login: logout*, May 2013.

[85] Evan Gilman: “The Discovery of Apache ZooKeeper’s Poison Packet,” *pagerduty.com*, May 7, 2015.

[86] Jonathan Stone and Craig Partridge: “When the CRC and TCP Checksum



Disagree,” at *ACM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication (SIGCOMM)*, August 2000. doi: 10.1145/347059.347561.

[87] Evan Jones: “How Both TCP and Ethernet Checksums Fail,” *evanjones.ca*, October 5, 2015.

[88] Cynthia Dwork, Nancy Lynch, and Larry Stockmeyer: “Consensus in the Presence of Partial Synchrony,” *Journal of the ACM*, volume 35, number 2, pages 288-323, April 1988. doi:10.1145/42282.42283.

[89] Peter Bailis and Ali Ghodsi: “Eventual Consistency Today: Limitations, Extensions, and Beyond,” *ACM Queue*, volume 11, number 3, pages 55-63, March 2013. doi:10.1145/2460276.2462076.

[90] Bowen Alpern and Fred B. Schneider: “Defining Liveness,” *Information Processing Letters*, volume 21, number 4, pages 181-185, October 1985. doi: 10.1016/0020-0190(85)90056-0.

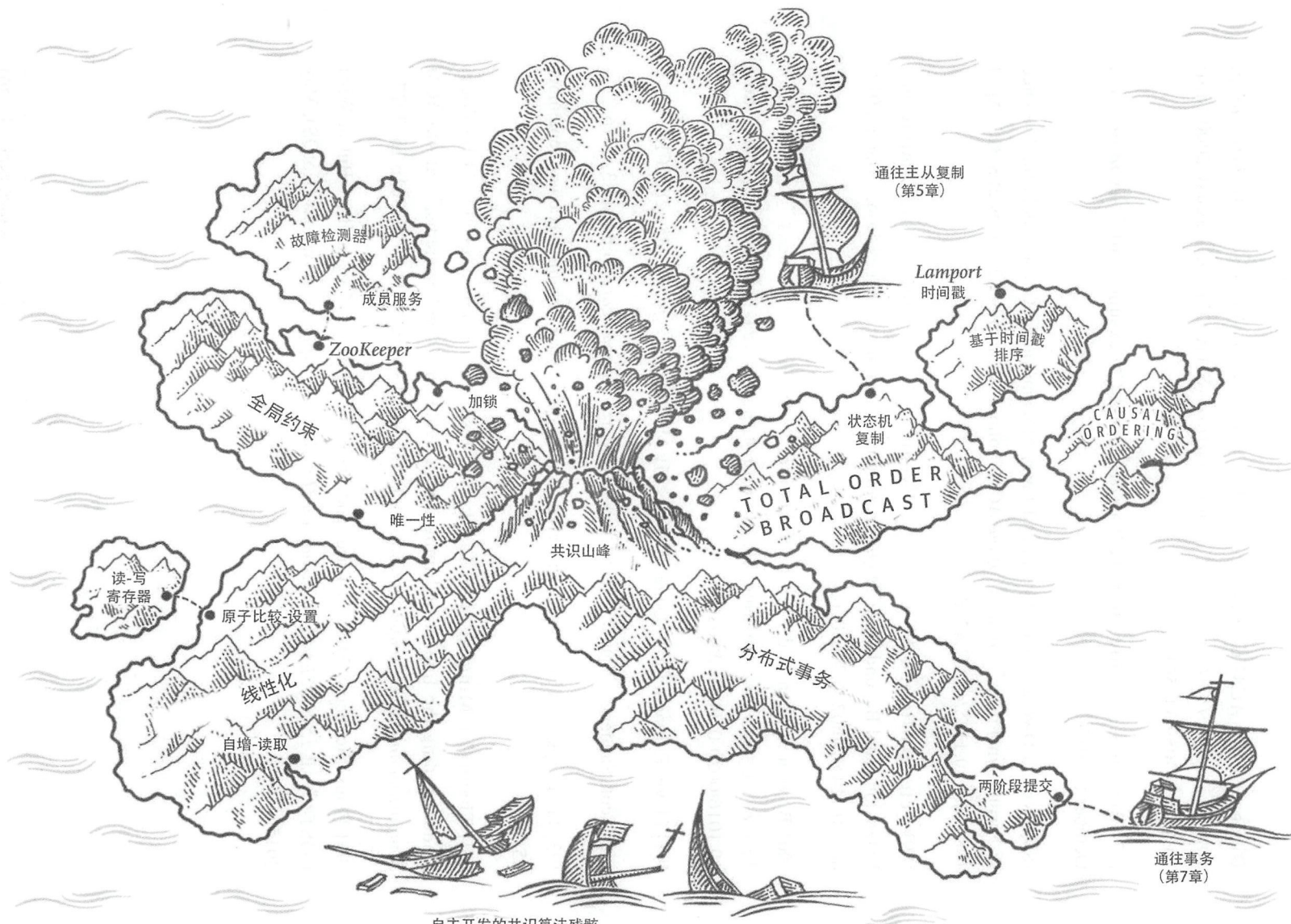
[91] Flavio P. Junqueira: “Dude, Where’s My Metadata?,” *fjp.me*, May 28, 2015.

[92] Scott Sanders: “January 28th Incident Report,” *github.com*, February 3, 2016.

[93] Jay Kreps: “A Few Notes on Kafka and Jepsen,” *blog.empathybox.com*, September 25, 2013.

[94] Thanh Do, Mingzhe Hao, Tanakorn Leesatapornwongsa, et al.: “Limplock: Understanding the Impact of Limpware on Scale-out Cloud Systems,” at *4th ACM Symposium on Cloud Computing (SoCC)*, October 2013. doi: 10.1145/2523616.2523627.

[95] Frank McSherry, Michael Isard, and Derek G. Murray: “Scalability! But at What COST?,” at *15th USENIX Workshop on Hot Topics in Operating Systems (HotOS)*, May 2015.



通往主从复制  
(第5章)

Lamport  
时间戳

基于时间戳  
排序

CAUSAL  
ORDERING

状态机  
复制

TOTAL ORDER  
BROADCAST

共识山峰

分布式事务

两阶段提交

通往事务  
(第7章)

自主开发的共识算法残骸

全局约束

加锁

唯一性

读-写  
寄存器

原子比较-设置

线性化

自增-读取

故障检测器

成员服务

ZooKeeper

# 一致性与共识

错误地活着，还是正确地挂掉？

——Jay Kreps，关于Kafka和Jepsen的若干注解

正如第8章所述，分布式系统存在太多可能出错的场景。而处理故障最简单的办法就是直接让整个服务停下来，然后向用户提示出错信息。但如果不能接受服务中止，就需要更加容错的解决方案，这样即使某些内部组件发生了故障，整个系统依然可以对外提供服务。

本章我们将讨论构建容错式分布式系统的相关算法和协议。这里假设第8章中所有的故障都可能发生，这包括网络数据包可能会丢失、顺序紊乱、重复发送或延迟，时钟也有一定偏差，节点可能发生暂停（例如由于垃圾回收）甚至随时崩溃。

为了构建容错系统，最好先建立一套通用的抽象机制和与之对应的技术保证，这样只需实现一次，其上的各种应用程序都可以安全地信赖底层的保证。这与第7章我们引入事务的道理相同：通过事务，应用程序可以假装没有崩溃（原子性），没有与其他人并发访问数据库（隔离性），且存储设备是完全可靠的（持久性）；总之，抽象的事务机制可以屏蔽系统内部很多复杂的问题，例如发生崩溃、边界条件、磁盘故障等，使得应用层轻松无忧。

现在继续沿着这个思路，尝试建立可以让分布式应用忽略内部各种问题的抽象机制。例如，分布式系统最重要的抽象之一就是共识：所有的节点就某一项提议达成一致。通过本章的介绍，最后你会发现面对各种网络故障和进程失效，可靠地达成共识是一件多么了不起的事情。



一旦解决了共识问题，就可以服务于应用层很多的目标需求。例如，对于一个主从复制的数据库，如果主节点发生失效，就需要切换到另一个节点，此时数据库节点可以采用共识算法来选举新的主节点。正如第5章“处理节点失效”所强调的，某一时刻必须只有一个主节点，所有的节点必须就此达成一致。如果有两个节点都自认为是主节点，就会发生脑裂，导致数据丢失。正确实现共识算法则可以避免此类问题。

本章我们将主要研究解决共识问题的相关算法。在此之前，首先我们会简要讨论分布式系统可提供的若干保证和抽象机制。

我们需要了解系统能力的边界，即哪些可行，哪些不可行。在什么情况下，系统可以容忍故障并继续工作；而其他情况，却无法保证。在理论证明和具体实现两方面，业界对此都有着非常深入的研究，在此本章即将为大家总体介绍各种系统边界情况。

事实上，几十年来分布式系统领域的研究人员一直在持续探索，积累了大量的研究资料，由于篇幅所限，我们无法一一触及。所以本章会省去相关的形式化模型与证明细节，更多的是给出较为直观的示例和解释。如有兴趣，本章后面的“参考文献”提供了更详细的信息。

## 一致性保证

在第5章“复制滞后问题”中，我们探讨了复制数据库时的计时问题。如果在同一时刻查询数据库的两个节点，则可能会看到不同的数据，这主要是因为写请求会在不同的时间点到达不同的节点。无论数据库采用何种复制方法（包括主从复制，多主节点复制或者无主节点复制），都无法完全避免这种不一致情况。

大多数多副本的数据库都至少提供了最终的一致性，这意味着如果停止更新数据库，并等待一段时间（长度未知）之后，最终所有读请求会返回相同的内容<sup>[1]</sup>。换句话说，不一致现象是暂时的，最终会达到一致（假设网络故障最终会被修复）。换言之，最终一致性意味着“收敛”，即预期所有的副本最终会收敛到相同的值<sup>[2]</sup>。

但是，这是一个非常弱的保证，它无法告诉我们系统何时会收敛。而在收敛之前，读请求可能会返回任何值甚至读失败<sup>[1]</sup>。例如，如果完成一笔更新操作之后立即读取，由于读取可能会路由到不同的副本，系统不保证一定读到刚刚写入的值（参阅第5章“读自己的写”）。

对于应用开发人员而言，最终一致性会带来很大的处理挑战，这与普通的单线程程序中变量读写行为大相径庭。对于后者，如果将对变量赋予某个值，接下来去读变量，

不可能发生读失败或者读不到刚赋的值。数据库表面上看起来像一个可以进行读写的变量，但事实上它内部有无比复杂的更多语义要求<sup>[3]</sup>。

当面对只提供了弱保证的数据库时，需要清醒地认清系统的局限性，切不可过于乐观。应用可能在大多数情况下都运行良好，但数据库内部可能已经发生了非常微妙的错误，只有当系统出现故障（例如网络中断）或高并发压力时，最终一致性的临界条件或者错误才会对外暴露出来，因而测试与发现错误变得非常困难。

因此本章将探索更强的一致性模型。不过，这也意味着更多的代价，例如性能降低或容错性差。尽管如此，更强的保证的好处是使上层应用逻辑更简单，更不容易出错。当了解、对比了多种不同的一致性模型之后，可以结合自身需求，从中选择最合适的。

分布式一致性模型与我们之前讨论过的多种事务隔离级别有相似之处<sup>[4,5]</sup>（参阅第7章“弱隔离级别”）。虽然存在某些重叠，但总体讲他们有着显著的区别：事务隔离主要是为了处理并发执行事务时的各种临界条件，而分布式一致性则主要是针对延迟和故障等问题来协调副本之间的状态。

本章内容涵盖非常广泛，但细究起来，这些内容之间存在着密切联系：

- 我们首先介绍线性化，这是最强的一致性模型，并考察其优缺点。
- 然后，我们将探讨分布式系统中事件顺序问题（参阅本章后面的“顺序保证”），特别是因果关系和全局顺序。
- 最后“分布式事务与共识”小节，我们将探索如何自动提交分布式事务，并最终解决共识问题。

## 可线性化

在最终一致性数据库中，同时查询两个不同的副本可能会得到两个不同的答案。这会使应用层感到困惑。如果数据库能够对上提供只有单个副本的假象，情况会不会大为简化呢？这样让每个客户端都拥有相同的数据视图，而不必担心复制滞后。

这就是可线性化<sup>[6]</sup>（也称为原子一致性<sup>[7]</sup>，强一致性等<sup>[8]</sup>）的思想。线性化的确切定义比较微妙，我们将稍后再详细探讨。其基本的想法是让一个系统看起来好像只有一个数据副本，且所有的操作都是原子的。有了这个保证，应用程序就不需要关心系统内部的多个副本。



在一个可线性化的系统中，一旦某个客户端成功提交写请求，所有客户端的读请求一定都能看到刚刚写入的值。这种看似单一副本的假象意味着它可以保证读取最近最新值，而不是过期的缓存。换句话说，可线性化是一种就近的保证。为了解释该想法，我们先来看一个非线性化系统的例子。

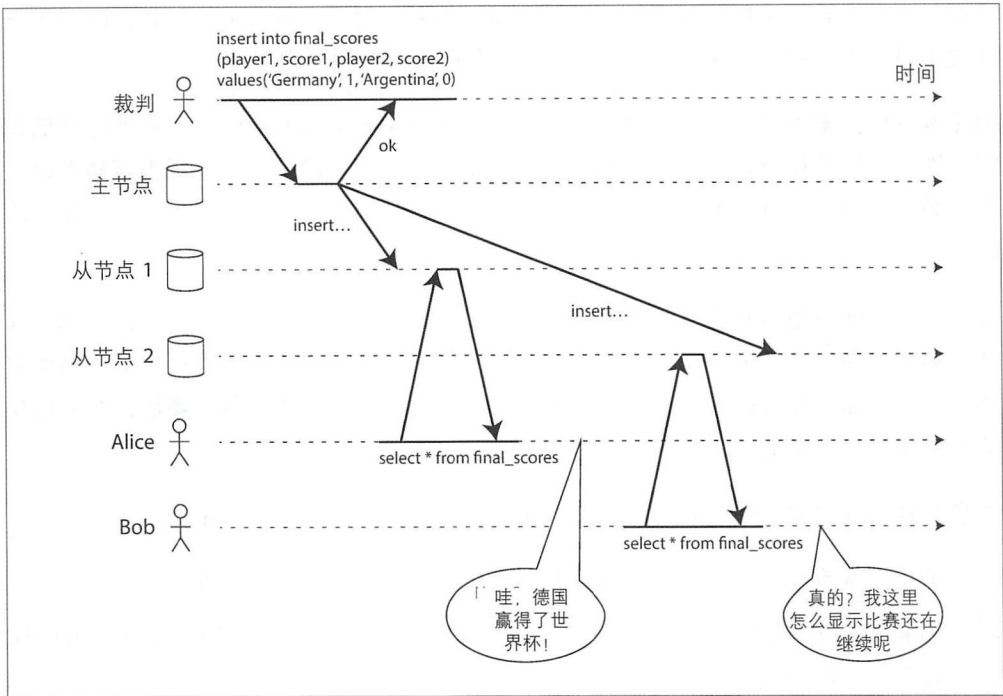


图9-1：一个非线性化的体育网站系统

图9-1是一个非线性化的体育网站<sup>[9]</sup>。Alice和Bob坐在同一个房间里各自观看自己的手机，焦急地等待2014年FIFA世界杯决赛的结果。刚刚宣布了最终比分之一后，Alice刷新了页面马上看到获胜者，然后兴奋地告诉Bob。于是Bob马上在自己的手机上刷新页面，但他的请求发向了某个落后的数据库副本，结果却显示比赛还在进行之中。

如果Alice和Bob几乎同时刷新页面得到两个不同的结果，他们也不清楚服务器端究竟何时接受、处理这些请求，因此可能并不会特别惊讶。然而，现在的情况却是Bob听到了Alice兴奋的比分之一后再单击刷新，他希望至少是看到刚刚Alice播报的最近比分之一，但却读到过期的结果，这就违背了线性化规则。





## 如何达到线性化?

可线性化背后的基本思想很简单：使系统看起来好像只有一个数据副本。然而，细究起来还有更多的含义。为了更好地理解可线性化，来看看更多的例子。

图9-2展示了三个客户端在线性化数据库中同时读写相同的主键 $x$ 。在分布式语义下， $x$ 被称为寄存器，例如，它可以是键-值存储中的一个键，关系数据库中的一行或文档数据库中的一个文档。

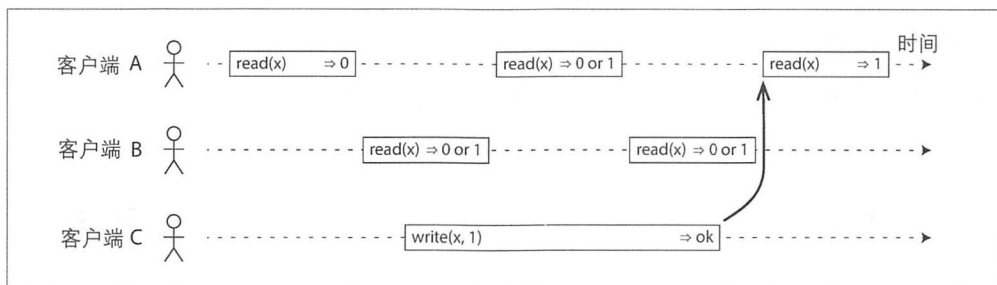


图9-2：读请求与写请求并发时，则可能读到旧值

为简单起见，图9-2仅展示了客户端所观察到请求内容，而不是数据库内部。每条线代表一个客户端请求，虚线的开始表示发送请求的时间，结尾则是收到响应的时间。由于网络延迟不确定，客户端并不清楚数据库具体何时处理请求，而只知道它是在发送之后、响应之前的某个中间时间点<sup>注1</sup>。

在这个例子中，寄存器有两类操作：

- $read(x) \Rightarrow v$  表示客户端读取 $x$ 的值，数据库返回了值 $v$ 。
- $write(x, v) \Rightarrow r$  表示客户端将 $x$ 设置为值 $v$ ，数据库返回处理结果 $r$ （可能表示处理成功或者发生了失败）。

在图9-2中， $x$ 的初始值为0，客户端C提交写请求将其设置为1。同时，客户端A和B在反复轮询数据库以读取最新值。A和B可能会分别读到什么样的返回值呢？

- 客户端A的第一个读取操作在写入开始之前已完成，因此返回的是旧值0。
- 客户端A的最后一次读操作是在写操作完成之后才开始的，如果数据库是可线性

注1： 图中的一个细节是，假设存在一个全局时钟（用水平轴表示）。即使实际的系统通常没有精确的时钟（参阅第7章“不可靠的时钟”），但这个假设并不妨碍此处对算法的分析：我们可以假设存在精确的全局时钟，只要算法并不实际去访问它<sup>[47]</sup>。



- 化的，它肯定会返回新值1。道理很简单，执行写操作肯定是在写请求发送之后并且响应之前；执行读操作同理。如果写入结束后开始读取，那么读一定发生在执行写之后，所以看到的一定是写入的新值。
- 与写操作有时间重叠的任何读取操作则可能返回0或者1，这是因为读写之间存在并发，无法确切知道在执行读取时，写入是否已经生效。

然而，这还没有精确描述线性化：如果与写并发的读操作可能返回旧值或新值，那么在这个过程中，不同的读客户端会看到旧值和新值之间来回跳变的情况。这肯定不符合我们所期望的模拟“单一数据副本”<sup>注2</sup>。

为使系统可线性化，我们需要添加一个重要的约束，如图9-3所示。

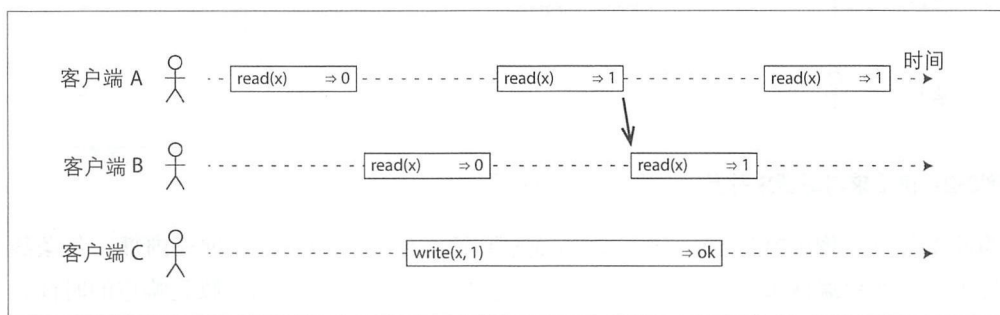


图9-3：一旦某个读操作返回了新值，之后所有的读（包括相同或不同的客户端）都必须返回新值

在一个可线性化的系统中，在写操作的开始与结束之间必定存在某个时间点，x的值发生了从0到1的跳变。如果某个客户端的读取返回了新值1，即使写操作尚未提交，那么所有后续的读取也必须全部返回新值。

图9-3中的箭头表示时序依赖关系。客户端A首先读到新值1，在A读取返回之后，B开始读取，由于B的读取严格在A的读取之后发生，因此即使C的写入仍在进行之中，也必须返回1。这和图9-1中Alice和Bob的情况类似，在Alice读取新值之后，Bob也预期读取新值。

可以进一步细化时序图来可视化每步操作具体在哪个时间点生效，如图9-4所示<sup>[10]</sup>。

在图9-4中，除了读写之外，我们引入了第三种类型的操作：

- $\text{cas}(x, v_{\text{old}}, v_{\text{new}}) \Rightarrow r$  表示一个原子比较-设置操作（compare-and-set, CAS）（参

注2：这种读写并发时，读可能返回旧值或新值的寄存器又被称为常规寄存器 [7,25]。



阅第7章“原子比较和设置”)。如果寄存器 $x$ 当前值等于 $v_{old}$ ,则将其原子设置为 $v_{new}$ ;否则保留现有 $x$ 值不变,然后返回错误, $r$ 是返回值,表示成功或者失败。

图9-4中的每个操作都有一条竖线,表示可能的执行时间点。这些标记以前后关系依次连接起来,最终的结果必须是一个有效的寄存器读写顺序,即每个读操作须返回最近写操作所设置的值。

可线性化要求,如果连接这些标记的竖线,它们必须总是按时间箭头(从左到右)向前移动,而不能向后移动。这个要求确保了之前所讨论的就近性保证:一旦新值被写入或读取,所有后续的读都看到的是最新的值,直到被再次覆盖。

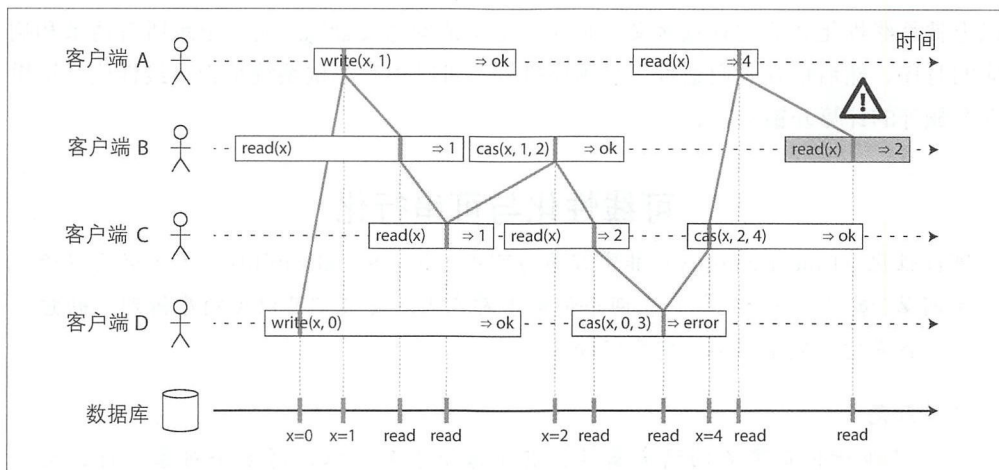


图9-4: 细化读、写入操作的生效时间, 客户端B的最后读取不满足可线性化

图9-4中有一些有趣的细节值得仔细分析:

- 客户端B首先发送读 $x$ 的请求, 接下来客户端D发送请求将 $x$ 置为0, 紧接着客户端A又发送请求将 $x$ 置为1, 而最终返回给B的值为1 (A所写入的值)。这是可能的, 它意味着数据库执行的顺序是: 首先处理D的写入0, 然后是A的写入1, 最后是B的读取。虽然这并不是请求发送的顺序, 但考虑到请求并发以及网络延迟等情况, 例如或许B的读请求网络延迟更大, 导致在两次写执行之后才到达数据库, 因此这是一个合法的可接受的处理顺序。
- 客户端A在收到数据库写响应之前, 客户端B即读到了值1, 这表明写入已成功。这是可能的, 但它并不代表执行读发生在执行写之前, 只是意味着很可能由于网络延迟而耽搁了客户端A接受响应。





- 模型没有假定事务间的隔离，即另一个并发客户端可能随时会修改值。例如，C首先读取到1，然后读到2，原因是两次读取之间值被客户端B修改了。我们可以使用原子比较和设置（cas）操作来检查值是否被其他并发客户端修改，例如客户端B和C的cas请求成功，但是D的cas操作失败（因为数据库处理时，x的值已不再为0）。
- 客户B的最后一次读取（阴影的方框）不满足线性化。该操作与C的cas写操作同时发生，后者将x从2更新为4。在没有其他请求时，B读取可以返回2。但是在B读取开始之前，客户端A已经读取了新值4，所以不允许B读到比A更老的值。这一点与图9-1中的Alice和Bob的情况是一样的。

以上就是线性化背后的直觉含义。正式的定义请参考文献<sup>[6]</sup>。通过记录所有请求和响应的时序，然后检查它们是否可以顺序排列，可以用来测试系统是否可线性化（这里存在额外的计算开销）<sup>[11]</sup>。

## 可线性化与可串行化

可线性化（Linearizability）非常容易与可串行化（Serializability）发生混淆（请参阅第7章“可串行化”），两个词似乎都在表达类似“可以按顺序排列”的意思。但是它们完全不同，需要仔细区分：

### 可串行化

可串行化是事务的隔离属性，其中每个事务可以读写多个对象（行，文档，记录等，请参阅第7章“单对象与多对象事务操作”）。它用来确保事务执行的结果与串行执行（即每次执行一个事务）的结果完全相同，即使串行执行的顺序可能与事务实际执行顺序不同<sup>[12]</sup>。

### 可线性化

可线性化是读写寄存器（单个对象）的最新值保证。它并不要求将操作组合到事务中，因此无法避免写倾斜等问题（请参阅第7章“写倾斜与幻读”），除非采取其他额外措施（如实现实体化冲突，参阅第7章“实体化冲突”）。

数据库可以同时支持可串行化与线性化，这种组合又被称为严格的可串行化或者强的单副本可串行化（strong one-copy serializability, strong-1SR）<sup>[4,13]</sup>。基于两阶段加锁（参阅第7章“两阶段加锁”）或者实际以串行执行（参阅第7章“实际串行执行”）都是典型的可线性化。



但是，可串行化的快照隔离（参阅第7章“可串行化的快照隔离”）则不是线性化的：按照设计，它可以从不一致性快照中读取，以避免读、写之间的竞争。一致性快照的要点在于它里面不包括快照点创建时刻之后的写入数据，因此从快照读取肯定不满足线性化。

## 线性化的依赖条件

那什么情况下应该使用线性化呢？上面足球比赛比分的例子只是个最简单的情况，结果存在几秒的延迟通常不会造成实质的伤害。然而，在有些场景下，线性化对于保证系统正确工作至关重要。

### 加锁与主节点选举

主从复制的系统需要确保有且只有一个主节点，否则会产生脑裂。选举新的主节点常见的方法是使用锁：即每个启动的节点都试图获得锁，其中只有一个可以成功即成为主节点<sup>[14]</sup>。不管锁具体如何实现，它必须满足可线性化：所有节点都必须同意哪个节点持有锁，否则就会出现问題。

提供协调者服务的系统如Apache ZooKeeper<sup>[15]</sup>和etcd<sup>[16]</sup>等通常用来实现分布式锁和主节点选举。它们都使用了支持容错的共识算法确保可线性化（本章后面的“支持容错的共识”会详细讨论这些算法）<sup>注3</sup>。虽然目前也有像Apache Curator<sup>[17]</sup>这样的上层库，在ZooKeeper之上提供了更高级别的接口以方便使用，但正确实现加锁和选举其实还有很多重要的设计细节（参阅第8章“主节点与锁”）。归根结底，线性化存储服务是所有这些协调服务的基础。

在一些分布式数据库如Oracle Real Application Clusters (RAC)<sup>[18]</sup>，分布式锁有更细粒度的实现：RAC为每个磁盘页面均设置一把锁，多个节点因此可以并发地共享访问存储系统。这些可线性化的锁处于事务执行的关键路径上，出于性能考虑，RAC部署时通常都要求专用的集群互连网络来连接数据库节点。

### 约束与唯一性保证

唯一性约束在数据库中很常见。例如，用户名或电子邮件地址必须唯一标识一个用

注3：严格来说，ZooKeeper和etcd提供可线性化的写操作，默认情况下，由于读操作可以在任何一个副本上执行，因此可能会读到过期值。如果需要，也可以激活线性化读，在etcd中称之为法定票数读取<sup>[16]</sup>，在ZooKeeper中，则需要读取之前调用sync()<sup>[15]</sup>。具体请参阅本章后面的“使用全序广播实现线性化存储”。





户，文件存储服务中两个文件不能具有相同的路径和文件名。如果要在写入数据时强制执行这些约束（例如，如果两个人试图同时创建具有相同名称的用户或文件，其中一个必须返回错误），则也需要线性化。

这种情况本质上与加锁非常类似：用户注册等同于试图对用户名进行加锁操作。该操作也类似于原子比较和设置：如果当前用户名尚未被使用，就设置用户名与客户ID进行关联。

其他类似约束包括银行账户余额不应出现负值，或者避免出售库存里已经没有的商品，或者不能同时预定航班或者剧院的相同的座位。这样的约束条件都要求所有节点就某个最新值达成一致（例如账户余额，库存水平，座位占用率）。

当然在某些实际场合中，有时可以放宽这些限制（例如，如果航班发生超额预订，可以将客户转移到其他的航班并提供必要的补偿）。在这种情况下，或许不需要线性化，我们将在第12章“时效性与完整性”讨论这些松散的约束条件。

然而，硬性的唯一性约束，常见如关系型数据库中主键的约束，则需要线性化保证。其他如外键或属性约束，则并不要求一定线性化<sup>[19]</sup>。

## 跨通道的时间依赖

请注意图9-1中的一个细节：如果Alice没有高呼比分，Bob可能就不会知道他的查询结果是过期的。或许他会在几秒之后再次刷新页面，然后看到最终的比方。线性化违例之所以被注意到，是因为系统中存在其他的通信渠道（例如，Alice对Bob发出的声音来传递信息）。

计算机系统也会出现类似的情况。例如，用户可以上传照片到某网站，有一个后台进程将照片调整为更低的分辨率（即缩略图）以方便更快下载。该网站架构和数据流如图9-5所示。

这里需要明确通知图像调整模块来调整哪些图片，系统采用了消息队列将此命令从Web服务器发送到调整器。因为大多数消息队列系统并不适合大数据流，而考虑到照片的大小可能到数兆字节，因此Web服务器并不会把照片直接放在队列中。相反，照片会先写入文件存储服务，当写入完成后，把调整的命令放入队列。

如果文件存储服务是可线性化的，那么系统应该可以正常工作。否则，这里就会引入竞争条件险：消息队列（图9-5中的步骤3和步骤4）可能比存储服务内部的复制执行更快。在这种情况下，当调整模块在读取图像（步骤5）时，可能会看到图像的某个







旧版本，或者根本读不到任何内容。如果它碰巧读到了旧版本的图像并进行处理，会导致文件存储中的全尺寸图片与调整之后图片出现永久的不一致。

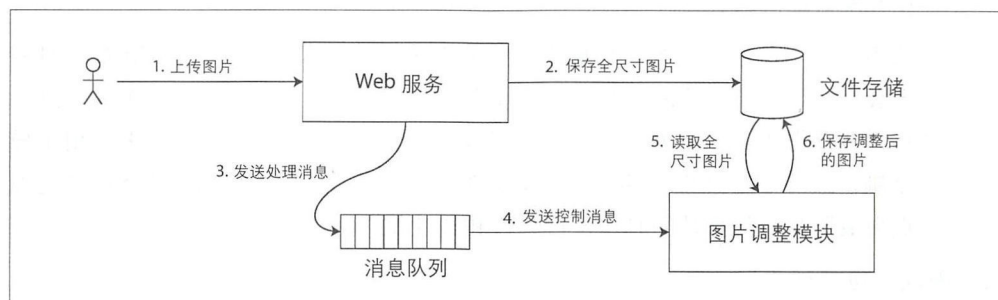


图9-5：Web服务器和图像调整器通过文件存储和消息队列进行通信，存在边界条件的可能

之所以出现这个问题是因为Web服务器和调整模块之间存在两个不同的通信通道：文件存储器和消息队列。如果没有线性化的就近性保证，这两个通道之间存在竞争条件。这种情况类似于图9-1：网页数据库有数据复制通道，而在计算机系统之外，Alice和Bob之间还有口耳相传的通道。

线性化并非避免这种竞争的唯一方法，但却是最容易理解的。如果可以控制某一个通信通道（例如消息队列，但注意不适合Alice和Bob的例子，因为后者并非计算机系统），可以尝试第5章的“读自己的写”方法，但会引入额外的复杂性。

## 实现线性化系统

现在我们已经看了几个线性化的例子，接下来考虑如何实现这样可线性化的系统。

由于线性化本质上意味着“表现得好像只有一个数据副本，且其上的所有操作都是原子的”，所以最简单的方案自然是只用一个数据副本。但显然，该方法无法容错：如果仅有的副本所在的节点发生故障，就会导致数据丢失，或者至少在重启之前都无法访问服务。

系统容错最常见的方法就是采用复制机制。我们再来回顾一下第5章所介绍的多种复制方案，看看哪些满足可线性化：

### 主从复制（部分支持可线性化）

在主从复制的系统中（参阅第5章的“主节点与从节点”），只有主节点承担数据写入，从节点则在各自节点上维护数据的备份副本。如果从主节点或者同步更



新的从节点上读取，则可以满足线性化<sup>注4</sup>。但并非每个主从复制的具体数据库实例都是可线性化的，主要是因为它们可能采用了快照隔离的设计，或者实现时存在并发方面的bug<sup>[10]</sup>。

而从主节点上读取的前提是你确定知道哪个节点是主节点。正如在第8章“真相由多数决定”中所讨论的，某节点可能自认为是主节点，但事实并非如此，这个“自以为是”的主节点如果对外提供服务，就会违反线性化<sup>[20]</sup>。如果使用了异步复制，故障切换过程中甚至可能会丢失一些已提交的写入（参阅第5章“处理节点失效”），结果是同时违反持久性和线性化。

### 共识算法（可线性化）

我们本章稍后即将讨论的一些共识算法，与主从复制机制相似。不过共识协议通常内置一些措施来防止裂脑和过期的副本。正是由于这些专门的设计，共识算法可以安全地实现线性化存储，这些系统包括ZooKeeper<sup>[21]</sup>和etcd<sup>[22]</sup>等。

### 多主复制（不可线性化）

具有多主节点复制的系统通常无法线性化的，主要由于它们同时在多个节点上执行并发写入，并将数据异步复制到其他节点。因此它们可能会产生冲突的写入，需要额外的解决方案（参阅第5章的“处理写冲突”）。这类冲突其实正是多副本所引入的结果。

### 无主复制（可能不可线性化）

对于无主节点复制的系统（即Dynamo风格，参阅第5章的“无主节点复制”），有些人认为只要配置法定读取和写入满足（ $w + r > n$ ）就可以获得“强一致性”。但这完全取决于具体的quorum的配置，以及如何定义强一致性，它可能并不保证线性化。

例如基于墙上时钟（包括Cassandra，参阅第8章“依赖于同步的时钟”）的“最后写入获胜”冲突解决方法几乎肯定是非线性化，因为这种时间戳无法保证与实际事件顺序一致（例如由于时钟偏移）。不规范的quorum（参阅第5章“宽松的quorum与数据回传”）也会破坏线性化。甚至即使是严格的quorum，正如之后即将介绍的，也会发生违背线性化的情况。

## 线性化与quorum

直觉上，对于Dynamo风格的复制模型，如果读写遵从了严格quorum，应该是可线性

---

注4： 如果对主从复制的数据库进行了分区（分片），然后每个分区上只有唯一的主节点，并不会影响可线性化。跨分区事务则是另一回事（请参阅本章后面的“分布式事务与共识”）。



化的。然而如果遭遇不确定的网络延迟，就会出现竞争条件，如图9-6所示。

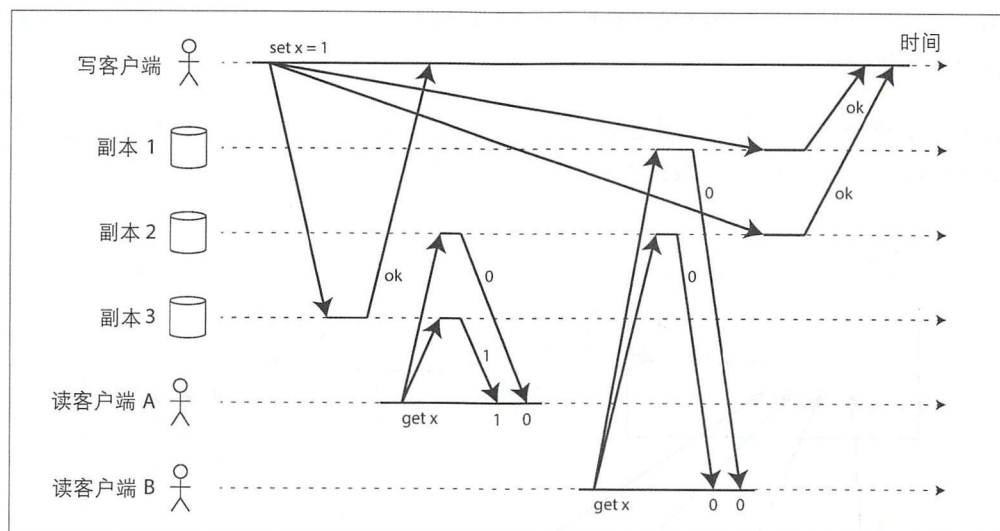


图9-6：尽管使用了严格的quorum，仍然不满足线性化

图9-6中， $x$ 的初始值为0，写客户端向所有三个副本（ $n = 3$ ， $w = 3$ ）发送写请求将 $x$ 更新为1。与此同时，客户端A从两个节点（ $r = 2$ ）读取数据，然后在其中一个节点上看到了新值1。与此同时，客户端B从两个节点的读取，两者在都返回了旧值0。

我们发现它虽然满足了仲裁条件（ $w + r > n$ ），但很明显这不是线性化的：B的请求在A的请求完成之后才开始，A返回了新值，但B却得到了旧值。这又类似图9-1中Alice和Bob的情况。

有趣的是，可以使Dynamo风格的复制系统以牺牲性能为代价来满足线性化：读操作在返回结果给应用之前，必须同步执行读修复（参阅第5章“读修复与反熵”）；而写操作在发送结果之前，必须读取quorum节点以获取最新值<sup>[24,25]</sup>。然而，由于会显著降低性能，Riak并不支持同步读修复<sup>[26]</sup>；Cassandra确实会等待读修复完成<sup>[27]</sup>，但是，它使用了“最后写入获胜”冲突解决方案，当出现同一个主键的并发写入时，就会丧失线性化。

此外，这种方式只能实现线性化读、写操作，但无法支持线性化的“比较和设置”操作，后者需要共识算法的支持<sup>[28]</sup>。

总而言之，最安全的假定是类似Dynamo风格的无主复制系统无法保证线性化。





## 线性化的代价

由于有一部分复制方案能够保证线性化，而其他则无法保证，因此有必要更加深入地探讨线性化的优缺点。

在第五章我们已经讨论了不同复制方案各自适合的场景。例如，多主复制非常适合多数据中心（参阅第5章“多数据中心操作”）。图9-7给出了这样的部署例子。

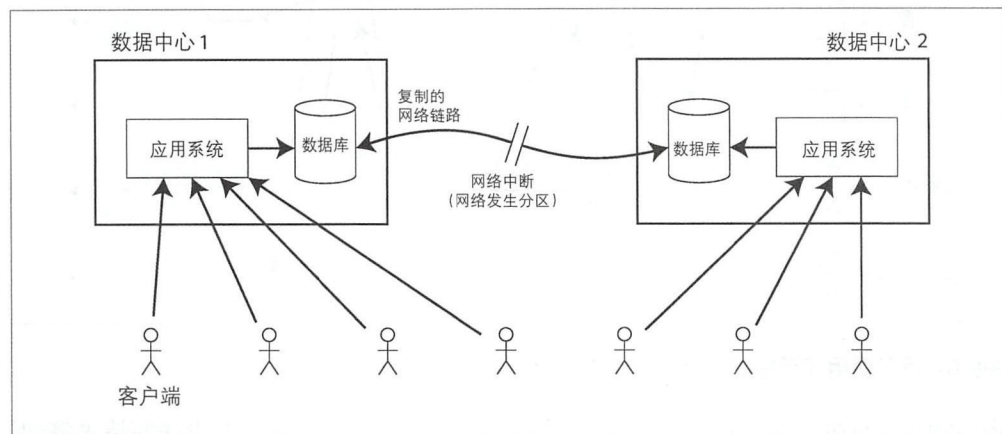


图9-7：网络中断迫使在可线性化与可用性之间做出选择

如果两个数据中心之间发生网络中断，会发生什么情况？我们假设每个数据中心内的网络工作正常，客户端可以到达就近的数据中心，但数据中心之间却无法互通。

基于多主复制的数据库，每个数据中心内都可以继续正常运行：由于从一个数据中心到另一个数据中心的复制是异步，期间发生的写操作都暂存在本地队列，等网络恢复之后再继续同步。

与之对比，如果是主从复制，则主节点肯定位于其中的某一个数据中心。所有写请求和线性化读取都必须发送给主节点，因此，对于那些连接到非主节点所在数据中心的客户端，读写请求都必须通过数据中心之间的网络，同步发送到主节点所在的数据中心。

因此，对于这样的主从复制系统，数据中心之间的网络一旦中断，连接到从数据中心的客户端无法再联系上主节点，也就无法完成任何数据库写入和线性化读取。从节点可以提供读服务，但内容可能是过期的（非线性化保证）。所以，如果应用程序要求线性化读写，则网络中断一定会违背这样的要求。

另一种情况，如果客户端可以直接连接到主节点所在的数据中心，则可以避免此问



题。否则，只能等到数据中心之间的网络恢复之后才能继续正常工作。

## CAP理论

不仅仅是主从复制和多主复制才有上面的问题，无论如何实现，任何可线性化的数据库都有这样问题；事实上，这个问题也不局限于多数据中心部署的情况，即使在一个数据中心内部，只要有不可靠的网络，都会发生违背线性化的风险。我们可以做以下的权衡考虑<sup>注5</sup>：

- 如果应用要求线性化，但由于网络方面的问题，某些副本与其他副本断开连接之后无法继续处理请求，就必须等待网络修复，或者直接返回错误。无论哪种方式，结果是服务不可用。
- 如果应用不要求线性化，那么断开连接之后，每个副本可独立处理请求例如写操作（多主复制）。此时，服务可用，但结果行为不符合线性化。

因此，不要求线性化的应用更能容忍网络故障。这种思路通常被称为CAP定理<sup>[29-32]</sup>，它由Eric Brewer于2000年正式命名，但基本想法可追溯到20世纪70年代，很多分布式数据库设计者当时都已经注意到这种现象<sup>[33-36]</sup>。

CAP最初是作为一个经验法则而提出的，并没有准确的定义，目的也只是帮助大家深入探讨数据库设计的权衡之道。当时，许多分布式数据库仍热衷于在集中共享的存储集群上提供可线性化的语义<sup>[18]</sup>，而CAP则事实上在鼓励大家去探索无共享系统，后者更适合于大规模的Web服务<sup>[37]</sup>，拥有更广阔的发展前景。CAP确实值得称赞，正是由于这种思路的转变，我们才见证了2000年以来新型数据库技术爆炸式的增长（可以概括为NoSQL系统）。

### CAP理论是否有用？

CAP有时也代表一致性，可用性，分区容错性，系统只能支持其中两个特性。不过，这种理解存在误导性<sup>[32]</sup>，网络分区是一种故障，不管喜欢还是不喜欢，它都可能发生，所以无法选择或逃避分区的问题。

在网络正常的时候，系统可以同时保证一致性（线性化）和可用性。而一旦发生了网络故障，必须要么选择线性（一致性），要么可用性。因此，更准确的称呼

注5：这两种选择有时分别称为CP（网络分区情况下选择一致性而丧失可用性）和AP（网络分区下可用但是不一致）。然而这种分类方案存在一些缺陷<sup>[9]</sup>，还是建议避免使用。



应该是“网络分区情况下，选择一致还是可用”<sup>[39]</sup>。高可靠的网络会帮助减少发生的概率，但无法做到彻底避免。

有必要指出，在CAP的诸多讨论中，术语可用性存在争议，其形式化定理中的可用性<sup>[30]</sup>与通常意义上的理解有些差别<sup>[40]</sup>。许多所谓的“高可用性”（容错）系统实际上并不符合CAP对可用性的特殊定义。总之，围绕着CAP有太多的误解与困扰，最后反而无法帮助我们更好地理解系统，所以本人建议最好避免使用CAP。

正式定义的CAP定理<sup>[30]</sup>范围很窄，它只考虑了一种一致性模型（即线性化）和一种故障（网络分区<sup>注6</sup>，节点仍处于活动状态但相互断开），而没有考虑网络延迟、节点失败或其他需要折中的情况。因此，尽管CAP在历史上具有重大的影响力，但对于一个具体的系统设计来说，它可能没有太大的实际价值<sup>[9,40]</sup>。

分布式系统中还有很多有趣的研究结果<sup>[41]</sup>，目前CAP已被更精确的研究成果所取代<sup>[2,42]</sup>，所以它现在更多的是代表历史上曾经的一个关注热点而已。

## 可线性化与网络延迟

虽然线性化是个很有用的保证，但实际上很少有系统真正满足线性化。例如，现代多核CPU上的内存甚至就是非线性化<sup>[43]</sup>：如果某个CPU核上运行的线程修改一个内存地址，紧接着另一个CPU核上的线程尝试读取，则系统无法保证可以读到刚刚写入的值，除非使用了内存屏障或fence指令<sup>[44]</sup>。

出现这种现象的原因是每个CPU核都有自己独立的cache和寄存器。内存访问首先进入cache系统，所有修改默认会异步地刷新到主存。由于访问cache比访问主存要快得多<sup>[45]</sup>，所以这样的异步刷新特性对于现代CPU的性能至关重要。但是，这就导致了多个数据副本（一个在主存，另外几个在不同级别的cache中），而副本更新是异步方式，无法保证线性化。

为什么这样呢？首先，CAP理论不适用于当今的多核-内存一致性模型：在计算机内部，我们通常假设通信是可靠的，例如我们不会假定一个CPU核在与其他核断开之后还能安然工作。之所以放弃线性化的原因就是性能，而不是为了容错。

注6：正如第8章“实践中的网络故障”中介绍的，本书曾使用了术语分区来表示将大数据集分解为较小的数据集（类似分片，参阅第6章）。而此处CAP的字母P虽然是同一个字母，但所表示的网络分区其实是一种特殊的网络故障，由此前后带来的混淆实在无法避免。





许多分布式数据库也是类似，它们选择不支持线性化是为了提高性能，而不是为了保住容错特性<sup>[46]</sup>。无论是否发生了网络故障，线性化对性能的影响都是巨大的。

那我们是否能找到一个更有效的线性化实现方案呢？目前看来答案是否定的。Attiya和Welch<sup>[47]</sup>证明如果想要满足线性化，那么读、写请求的响应时间至少要与网络中延迟成正比。考虑到多数计算机网络高度不确定的网络延迟（参阅第7章“超时与无限延迟”），线性化读写的性能势必非常差。虽然没有足够快的线性化算法，但弱一致性模型的性能则快得多，这种取舍对于延迟敏感的系统非常重要。在第12章，我们将讨论一些避免线性化但又可以确保正确性的方法。

## 顺序保证

我们之前曾说过，线性化寄存器对外呈现的好像只有一份数据拷贝，而且每一个操作似乎都是原子性生效。这意味着操作是按照某种顺序执行，如图9-4所展示的执行顺序。

顺序是本书反复出现的主题，某种程度也表明它确实是一个非常重要的基本概念。让我们简要回顾一下本书讨论顺序时涉及的前后上下文：

- 在第5章，我们看到主从复制系统中主节点的主要作用是确定复制日志中的写入顺序，这样使从节点遵从相同的顺序执行写入。如果没有这样的唯一主节点，则可能由于并发操作而引发冲突（参阅第5章“处理写冲突”）。
- 第7章中讨论的可串行化则是确保事务的执行结果与按照某种顺序方式执行一样。实现方式可以是严格顺序执行，或者允许并发但需要相应的冲突解决方案（例如加锁或冲突-中止）。
- 第8章讨论了分布式系统的时间戳与时钟（参阅第8章“依赖于同步的时钟”），试图将顺序引入到无序的操作世界，例如确定两个写操作哪一个先发生。

事实证明，排序、可线性化与共识之间存在着某种深刻的联系。尽管这些概念听起来比本书的其他部分更加理论和抽象，但它对于理解系统能做什么和不能做什么非常有帮助。

## 顺序与因果关系

之所以反复出现“顺序”问题，其中的一个原因是它有助于保持因果关系。我们已经有好几个这样的例子来反复说明因果关系的重要性：

- 在第5章“一致前缀读”（见图5-5）中，我们看到这样的例子，一个对话的观察者首先看到了问题的答案接着才是问题本身。这违背了我们对因果的直觉认识，因而感觉困惑：问题之所以被回答，一定是有问题在先，意味着回答的人一定是先看到了问题（当然，假定他们没有精神问题，也无法预知未来）。此时，问题与回答之间存在因果关系。
- 图5-9类似，三个主节点之间进行数据复制，由于网络延迟，一些写操作会覆盖其他的写入。从某个副本的角度来看，好像是发生了一个对不存在数据行的更新。这里的因果意味着首先必须先创建数据行，然后才能去更新。
- 在第5章“检测并发写”，如果有两个操作A和B，则它们之间一共有三种可能性：A发生在B之前，B发生在A之前，或者A和B并发。这种“A发生在B之前”其实是因果关系的另一种表示。如果A发生在B之前，意味着B可能已经知道了A，或者基于A的操作，或者依赖于A。如果A和B是并发关系，则它们之间不存在因果关系；换言之，互相不知道对方。
- 在事务的快照隔离上下文中（参阅第7章“快照隔离与可重复读”），事务是从一致性快照读取。这里的“一致性”又是什么含义呢？这意味着与因果关系一致：如果快照中包含了答案，则它也必须包含所提的问题<sup>[48]</sup>。这样才能确保在某个时间点观察数据库时符合因果关系：快照创建时刻点之前的所有数据都要可见，但此后发生的事件则不可见。读倾斜（例如图7-6的不可重复读取）则违反了因果关系因而读到了本不可见的数据库。
- 事务之间写倾斜的例子（参阅第7章“写倾斜与幻读”）也说明了因果关系。在图7-8中，Alice申请调班成功是因为事务以为Bob仍在值班，反之亦然。在这种情况下，调班动作的因果关系取决于当前是谁在值班。可序列化的快照隔离（参阅第7章“可串行化的快照隔离”）主要通过跟踪事务之间的因果依赖关系从而达到检测写倾斜的目的。
- 在Alice和Bob一起看足球的例子中（见图9-1），当听到Alice惊呼比分之后，Bob从服务器却得到过期的结果违背了因果关系。Alice之所以高呼是因为她先看到了最新的比分，因而道理上讲，Bob也应该看到和Alice一样的比分。另一个类似例子则是通过两个通道调整图片大小。

因果关系对所发生的事件施加了某种排序：发送消息先于收到消息；问题出现在答案之前等，或者就像在现实生活中一样，一件事情会导致另一件事情：一个节点根据读取的数据做出决定，然后写入结果，另一个节点读取写入的结果之后再写入新的内容，诸如此类。这些因果关系的依赖链条定义了系统中的因果顺序，即某件事应该发生另一件事情之前。



如果系统服从因果关系所规定的顺序，我们称之为因果一致性。例如，快照隔离提供了因果一致性：当从数据库中读数据时，如果查询到了某些数据，也一定能看到触发该数据的前序事件（假设期间没有发生删除操作）。

## 因果顺序并非全序

全序关系支持任何两个元素之间进行比较，即对于任意两个元素，总是可以指出哪个更大，哪个更小。例如，自然数符合全序关系，随便给出两个数字比如5和13，都可以进行比较。

但是，有些集合并不符合全序，例如集合 $\{a, b\}$ 大于集合 $\{b, c\}$ 么？因为它们都不是对方的子集，所以无法直接比较它们。我们称之为不可比较，数学集合只能是偏序。某些情况下，一个集合可以包含另一个（如果该集合包含另一个集合的所有元素），否则则无法比较。

全序和偏序的差异也会体现在不同的数据库一致性模型中：

### 可线性化

在一个可线性化的系统中，存在全序操作关系。系统的行为就好像只有一个数据副本，且每个操作都是原子的，这意味着对于任何两个操作，我们总是可以指出哪个操作在先。这种全序排列序如图9-4中的时间线所示。

### 因果关系

如果两个操作都没有发生在对方之前，那么这两个操作是并发关系（参阅第5章“Happens-before关系与并发”）。换言之，如果两个事件是因果关系（一个发生在另一个之前），那么这两个事件可以被排序；而并发的操作则无法排序比较。这表明因果关系至少可以定义为偏序，而非全序。

因此，根据这个定义，在可线性化数据存储中不存在并发操作，一定有一个时间线将所有操作都全序执行。可能存在多个请求处于等待处理的状态，但是数据存储保证了在特定的时间点执行特定的操作，所以是单个时间轴，单个数据副本，没有并发。

并发意味着时间线会出现分支和合并，而不同分支上的操作无法直接比较。第5章中我们给出了这种例子，如图5-14并非一条直线式的全序，而是多个不同的操作同时进行。图中的箭头表明这只是因果关系，即部分操作之间的偏序。

如果熟悉像Git这样的分布式版本控制系统，那么它们的版本历史非常类似于因果关系图。通常情况下，提交会以直线形式呈现，但有时会产生分支（特别是多个人同时在一个项目上工作时），当同时有多个提交时就需要进行合并。





## 可线性化强于因果一致性

那么因果序和可线性化之间是什么关系呢？答案是可线性化一定意味着因果关系：任何可线性化的系统都将正确地保证因果关系<sup>[7]</sup>。特别是，如果系统存在多个通信通道（见图9-5中的消息队列和文件存储服务），可线性化确保了因果关系会自动全部保留，而不需要额外的工作（比如在不同组件之间的传递时间戳）。

可线性化可以确保因果性这一结论，使线性化系统更加简单易懂而富有吸引力。但是，正如在“线性化的代价”将要阐述的，线性化会显著降低性能和可用性，尤其是在严重网络延迟的情况下（例如多数据中心）。正因如此，一些分布式数据系统已经放弃了线性化，以换来更好的性能，但也存在可能无法正确工作的风险。

好消息是线性化并非是保证因果关系的唯一途径，还有其它方法使得系统可以满足因果一致性而免于线性化所带来的性能问题。事实上，因果一致性可以认为是，不会由于网络延迟而显著影响性能，又能对网络故障提供容错的最强的一致性模型<sup>[2,42]</sup>。

在许多情况下，许多看似需要线性化的系统实际上真正需要的是因果一致性，后者的实现可以高效很多。基于这样的观察，研究人员正在探索新的数据库来保证因果关系，其性能与可用性特征与最终一致性类似<sup>[49-51]</sup>。

由于这些研究还处于早期阶段，大多数还没有投入到生产系统，也存在其他一些挑战需要克服<sup>[52,53]</sup>。但是，毫无疑问，这是未来非常有希望的方向。

## 捕获因果依赖关系

我们此处不会深入探讨非线性化系统是如何保证因果一致性的细节，而只是简单地介绍一些关键思想。

为保持因果关系，需要知道哪个操作发生在前。这里只需偏序关系，或许并发操作会以任意顺序执行，但如果一个操作发生在另一个操作之前，那么每个副本都应该按照相同的顺序处理。因此，当某个副本在处理一个请求时，必须确保所有因果在前的请求都已完成处理；否则，后面的请求必须等待直到前序操作处理完毕。

为了确定请求的因果依赖关系，我们需要一些手段来描述系统中节点所知道的“知识”。如果节点在写入Y时已经看到X值，则X和Y可能是属于因果关系。这种分析使用了类似针对欺诈指控刑事调查中的推理方法：某个CEO在做出决定事项Y时是否已经知道了消息X？

确定请求的先后顺序与第五章“检测并发写”中所讨论的技巧类似。后者针对的是无



主复制中的因果关系，该场景需要去检测对同一个主键的并发写请求，从而避免更新丢失。因果一致性则要更进一步，它需要跟踪整个数据库请求的因果关系，而不仅仅是针对某个主键。版本向量技术可以推广为一种通用的解决方案<sup>[54]</sup>。

为了确定因果关系，数据库需要知道应用程序读取的是哪个版本的数据。这就是为什么在图5-13中先前读操作的版本号在提交时要传回到数据库。SSI的冲突检测也是类似想法，如第7章“可串行化的快照隔离”所介绍的：当事务提交时，数据库要检查事务曾经读取的数据版本现在是否仍是最新的。为此，数据库需要跟踪事务读取了哪些版本的数据。

## 序列号排序

虽然因果关系很重要，但实际上跟踪所有的因果关系不切实际。在许多应用程序中，客户端在写入之前会先读取大量数据，系统无法了解之后的写入究竟是依赖于全部读取内容，还是仅仅是其中一小部分。但很明显，显式跟踪所有已读数据意味着巨大的运行开销。

这里还有一个更好的方法：我们可以使用序列号或时间戳来排序事件。时间戳不一定来自墙上时钟（或者物理时钟，但正如第8章所讨论的，物理时钟存在很多问题）。它可以只是一个逻辑时钟，例如采用算法来产生一个数字序列用以识别操作，通常是递增的计数器。

这样的序列号或时间戳非常紧凑（只有几字节的大小），但它们保证了全序关系。也就是说，每一个操作都有唯一的顺序号，并且总是可以通过比较来确定哪个更大（即操作发生在后）。

特别是，我们可以按照与因果关系一致的顺序来创建序列号<sup>注7</sup>：保证如果操作A发生在B之前，那么A一定在全序中出现在B之前（即A的序列号更小）。并行操作的序列可能是任意的。这样的全局排序可以捕获所有的因果信息，但也强加了比因果关系更为严格的顺序性。

在主从复制数据库中（参阅第5章“主节点与从节点”），复制日志定义了与因果关系一致的写操作全序关系。主节点可以简单地为每个操作递增某个计数器，从而为复制日志中的每个操作赋值一个单调递增的序列号。从节点按照复制日志出现的顺

---

注7：很容易产生与因果关系不一致的全序关系，但用处不大。例如，可以为每个操作生成随机UUID，并按照字典序比较UUID来定义全局顺序。这肯定是一个全序，但随机的UUID无法区分哪个操作发生在先，哪些操作是并发。



序来应用写操作，那结果一定满足因果一致性（虽然从节点的数据可能会滞后于主节点）。

## 非因果序列发生器

如果系统不存在这样唯一的主节点（例如可能是多主或者无主类型的数据库，或者数据库本身是分区的），如何产生序列号就不是那么简单了。实践中可以采用以下方法：

- 每个节点都独立产生自己的一组序列号。例如，如果有两个节点，则一个节点只生成奇数，而另一个节点只生成偶数。还可以在序列号中保留一些位用于嵌入所属节点的唯一标识符，确保不同的节点永远不会生成相同的序列号。
- 可以把墙上时间戳信息（物理时钟）附加到每个操作上<sup>[55]</sup>。时间戳可能是不连续的，但是只要它们有足够高的分辨率，就可以用来区分操作。“最后写获胜”的冲突解决方案也使用类似的方法（参阅第8章“时间戳与事件顺序”）。
- 可以预先分配序列号的区间范围。例如，节点A负责区间1~1000的序列号，节点B负责1001~2000。然后每个节点独立地从区间中分配序列号，当序列号出现紧张时就分配更多的区间。

上述三种思路都可行，相比于把所有请求全部压给唯一的主节点具有更好的扩展性。它们为每个操作生成一个唯一的、近似增加的序列号。不过，它们也都存在一个问题：所产生的序列号与因果关系并不严格一致。

所有这些序列号发生器都无法保证正确捕获跨节点操作的顺序，因而存在因果关系方面的问题：

- 每个节点可能有不同的处理速度，如每秒请求数。因此，某个节点产生偶数而另一个产生奇数，偶数的计数器产生速度可能落后于奇数的计数器，反之亦然。这样就无法准确地知道哪个操作在先。
- 物理时钟的时间戳会受到时钟偏移的影响，也可能导致与实际因果关系不一致。例如图8-3，后来发生的操作实际上被分配了一个较低的时间戳<sup>注8</sup>。
- 对于区间分配器，一个操作可能被赋予从1001~2000之间的某个序列号，而后发

---

注8：当然有可能使物理时钟时间戳与因果关系保持一致，如第8章“同步时钟与全局快照”，我们介绍了Google Spanner，它可以给出时钟偏差的范围，然后在写提交之前等待一定的间隔，以此来确保后发生的事件被赋予足够大的时间戳。不过大多数时钟不具备此类高精度条件。





生的操作则路由到另一个节点，拿到了某个1~1000之间的序列号，导致与因果序不一致。

## Lamport时间戳

刚才所描述的三个序列号发生器可能与因果关系存在不一致，但还有一个简单的方法可以产生与因果关系一致的序列号。它被称为兰伯特时间戳（Lamport timestamp），由Leslie Lamport于1978年提出<sup>[56]</sup>，该文献也是现代分布式系统领域被引用最多的经典论文之一。

图9-8给出了Lamport时间戳的示例。首先每个节点都有一个唯一的标识符，且每个节点都有一个计数器来记录各自已处理的请求总数。Lamport时间戳是一个值对（计数器，节点ID）。两个节点可能会有相同的计数器值，但时间戳中还包含节点ID信息，因此可以确保每个时间戳都是唯一的。

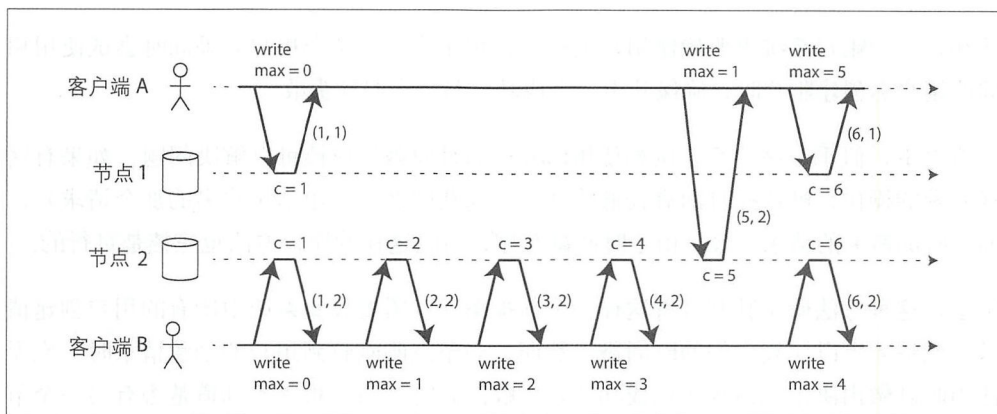


图9-8：Lamport时间戳可以保证全序与因果关系一致

Lamport时间戳与物理墙上时钟并不存在直接对应关系，但它可以保证全序：给定两个Lamport时间戳，计数器较大那个时间戳大；如计数器值正好相同，则节点ID越大，时间戳越大。

到目前为止，该思路与上一节所描述的奇数/偶数计数器并无本质不同。但是Lamport时间戳的核心亮点在于使它们与因果性保持一致，具体如下所示：每个节点以及每个客户端都跟踪迄今为止所见到的最大计数器值，并在每个请求中附带该最大计数器值。当节点收到某个请求（或者回复）时，如果发现请求内嵌的最大计数器值大于节点自身的计数器值，则它立即把自己的计数器修改为该最大值。

如图9-8所示，客户端A从节点2收到计数器值5，然后将最大值5发送到节点1。此时，





节点1的计数器仅为1，但是它立即向前跳到5，所以下一个操作将获得计数器值6。

只要把最大计数器值嵌入到每一个请求中，该方案可以确保Lamport时间戳与因果关系一致，而请求的因果依赖性一定会保证后发生的请求得到更大的时间戳。

Lamport时间戳有时会与版本向量发生混淆（第5章“检测并发写”中介绍了版本向量）。虽然存在一些相似之处，但它们的目的不同：版本向量用以区分两个操作是并发还是因果依赖，而Lamport时间戳则主要用于确保全序关系。即使Lamport时间戳与因果序一致，但根据其全序关系却无法区分两个操作属于并发关系，还是因果依赖关系。Lamport时间戳优于版本向量之处在于它更加紧凑和高效。

### 时间戳排序依然不够

虽然Lamport时间戳定义了与因果序一致的全序关系，但还不足以解决实际分布式系统中许多常见的问题。

例如，一个账户系统需要确保用户名唯一标识用户。即两个用户如果同时尝试使用相同的用户名创建账户时，确保其中一个成功，另一个必须失败。

乍看之下，似乎全序关系（例如使用Lamport时间戳）应该可以解决问题：如果有这样并发的操作，则选择时间戳较低的那个作为获胜者（先申请用户名的那个请求），而让时间戳大的请求失败。由于时间戳有序，所以这样的比较方法也应该是可行的。

但是，这种方法确定胜利者有这样一个前提条件：需要收集系统中所有的用户创建请求，然后才可以比较它们的时间戳。然而，当节点刚刚收到用户的创建请求时，它无法当时就做出决定该请求应该成功还是失败。此时，节点根本不知道是否有另一个节点在同时创建相同用户名（以及那个请求所附带的时间戳）。

而为了获得上述两点信息，系统就必须检查每个节点，询问它们在做什么<sup>[56]</sup>。如果万一某个节点出现故障或者由于网络问题而无法连接，那么方法就无法正常运转。显然这不是我们所期望的容错系统。

这里问题的关键是，只有在收集了所有的请求信息之后，才能清楚这些请求之间的全序关系。如果另一个节点执行了某些操作，但你无法知道那是什么，就无法构造出最终的请求序列。也许，来自该未知操作确实需要插入到全序集合中才能正确评估出下一步。

总而言之，为了实现像用户名唯一性约束这样的目标，仅仅对操作进行全序排列还是不够的，还需要知道这些操作是否发生、何时确定等。假如能够在创建用户名时，已





经确定知道了没有其他节点正在执行相同用户名的创建，你大可以直接安全返回创建成功。

要想知道什么时候全序关系已经确定就需要之后的“全序关系广播”。

## 全序关系广播

如果程序只运行在一个CPU核上，可以非常简单地定义出操作的全序关系，即在单核上执行的顺序。但是，在分布式系统中，让所有的节点就全序关系达成一致就面临巨大挑战。在前面一节中，我们讨论了按时间戳或序列号排序，发现它不如主从复制那么直接有效（如果使用时间戳排序来实现唯一性约束，会丧失容错性）。

如前所述，主从复制首先确定某一个节点作为主节点，然后在主节点上顺序执行操作。接下来的主要挑战在于，如何扩展系统的吞吐量使之突破单一主节点的限制，以及如何处理主节点失效时的故障切换（参阅第5章“处理节点失效”）。在分布式系统研究文献中，这些问题被称为全序关系广播或者原子广播<sup>[25,57,58]注9</sup>。



### 顺序保证的范围

每个分区只有一个主节点的数据库通常只会维护分区内的顺序，即它们不提供跨分区的一致性保证（例如，一致性快照，外键引用等）。跨分区的全序关系并非不可能，但需要非常多的额外工作<sup>[59]</sup>。

全序关系广播通常指节点之间交换消息的某种协议。下面是一个非正式的定义，它要求满足两个基本安全属性：

#### 可靠发送

没有消息丢失，如果消息发送到了某一个节点，则它一定要发送到所有节点。

#### 严格有序

消息总是以相同的顺序发送给每个节点。

即使节点或网络出现了故障，全序关系广播算法的正确实现也必须保证上述两条。当然，网络中断时是不可能发送成功的，但算法要继续重试，直到最终网络修复，消息发送成功（且必须以正确的顺序发送）。

注9：“原子广播”是个经典的术语，但却容易引起混淆，其中“原子”的含义与其他上下文含义并不一致。例如与ACID事务中的原子性没有任何关系，只是与原子操作（多线程编程上下文）或原子寄存器（线性化存储）有点间接的关系。







## 使用全序关系广播

像ZooKeeper和etcd这样的共识服务实际上就实现了全关系序广播。这也暗示了全序关系广播与共识之间有着密切联系，本章稍后会揭示这一点。

全序关系广播正是数据库复制所需要的：如果每条消息代表数据库写请求，并且每个副本都按相同的顺序处理这些写请求，那么所有副本可以保持一致（或许有些滞后）。该原则也被称为状态机复制<sup>[60]</sup>，我们将在第11章中详细介绍。

可以使用全序关系广播来实现可串行化事务。如第7章“实际串行执行”所述，如果每条消息表示一个确定性事务并且作为存储过程来执行，且每个节点都遵从相同的执行顺序，那么可以保证数据库各分区以及各副本之间的一致性<sup>[61]</sup>。

全序关系广播另一个要点是顺序在发送消息时已经确定，如果消息发送成功，节点不允许追溯地将某条消息插入到先前的某个位置上。这一点使得全序关系广播比基于时间戳排序要求更强。

理解全序关系广播的另一种方式是将其视为日志（如复制日志，事务日志或预写日志）。传递消息就像追加方式更新日志。由于所有节点必须以相同的顺序发送消息，因此所有节点都可以读取日志并看到相同的消息序列。

全序关系广播对于提供fencing令牌的锁服务也很有用（参阅第8章的“Fencing 令牌”）。每个获取锁的请求都作为消息附加到日志中，所有消息按照日志中的顺序依次编号。序列号还可以作为令牌，它符合单调递增要求。在ZooKeeper中，该序列号被称为zxid<sup>[15]</sup>。

## 采用全序关系广播实现线性化存储

如图9-4所示，在一个可线性化的系统中有全序操作集合。这是否意味着可线性化与全序关系广播是完全相同呢？不完全是，但两者之间有着密切的联系<sup>注10</sup>。

全序关系广播是基于异步模型：保证消息以固定的顺序可靠地发送，但是不保证消息何时发送成功（因此某个接收者可能明显落后于其他接收者）。而可线性化则强调就近性：读取时保证能够看到最新的写入值。

---

注10：通常意义上，实现线性化读写寄存器是一个“更容易”的问题。全序关系广播相当于共识<sup>[67]</sup>，在异步的崩溃-中止模型<sup>[68]</sup>中不存在确定性的解决方案，而线性化读写寄存器则可以在该模型下得以实现<sup>[23-25]</sup>。然而，如果进一步支持原子操作例如原子比较-设置，或者原子自增-读取，则它就升级等价于共识问题，具体请参考文献<sup>[28]</sup>。因此，共识与线性化的问题有着密切相关性。



如果有了全序关系广播，就可以在其上构建线性化的存储系统。例如，确保用户名唯一标识一个用户。

设想一下，对于每一个可能的用户名，都可以有一个带有原子比较-设置操作的线性化寄存器。每个寄存器初始值为空（表示尚未使用）。当用户创建一个用户名时，对该用户名的寄存器执行比较设置操作：仅当寄存器值为空时，将其设置为新的用户账号。如果多个用户试图同时获取相同的用户名，则只有一个原子比较-设置操作成功。

可以通过使用全序关系广播以追加日志的方式<sup>[62,63]</sup>来实现线性化的原子比较-设置操作，步骤如下所示：

1. 在日志中追加一条消息，并指明想要的用户名。
2. 读取日志，将其广播给所有节点，并等待回复<sup>注11</sup>。
3. 检查是否有任何消息声称该用户名已被占用。如果第一条这样的回复来自于当前节点，那么就成功获得该用户名，可以提交该获取声明（也许附加另一条消息到日志）并返回给客户端。反之，如果声称占用的第一条回复消息来自其他节点，则中止操作。

由于日志条目以相同的顺序发送到所有节点，而如果存在多个并发写入，则所有节点将首先决定哪个请求在先。选择第一个写请求作为获胜者，并中止其他请求，以确保所有节点同意一个写请求最终要么提交成功要么中止。类似的方法还可以用来在日志之上实现可串行化的多对象事务<sup>[62]</sup>。

虽然此过程可确保线性化写入，但它却无法保证线性化读取，即从异步日志更新的存储中读取数据时，可能是旧值。具体来说，这里只提供了顺序一致性<sup>[47,64]</sup>，有时也称为时间线一致性<sup>[65,66]</sup>，它弱于线性化保证。为了同时满足线性化读取，有以下几个方案：

- 可以采用追加的方式把读请求排序、广播，然后各个节点获取该日志，当本节点收到消息时才执行真正的读操作。消息在日志中的位置已经决定了读取发生的时间点。etcd的quorum读取和这个思路有相似之处<sup>[16]</sup>。
- 如果可以以线性化的方式获取当前最新日志中消息的位置，则查询位置，等待直

---

注11：如果不等待，而是入队之后立即确认成功，情况就会类似多核x86处理器上的内存一致性模型<sup>[43]</sup>。该模型不满足线性化也不满足顺序一致性。



到该位置之前的所有条目都已经发送给你，接下来再执行读取。这与ZooKeeper的sync()操作思路相同<sup>[15]</sup>。

- 可以从同步更新的副本上进行读取，这样确保总是读取最新值。这种技术可以用于链式复制<sup>[63]</sup>，具体参阅第5章的“复制研究”。

## 采用线性化存储实现全序关系广播

前面一节介绍了如何基于全序关系广播来构建线性化的原子比较-设置操作。我们也可以反过来，假定已有了线性化的存储，在其上构建全序关系广播。

最简单的方法是假设有一个线性化的寄存器来存储一个计数，然后使其支持原子自增-读取操作<sup>[28]</sup>或者原子比较-设置操作。

算法思路很简单：对于每个要通过全序关系广播的消息，原子递增并读取该线性化的计数，然后将其作为序列号附加到消息中。接下来，将消息广播到所有节点（如果发生丢失，则重新发送），而接受者也严格按照序列化来发送回复消息。

请注意，与Lamport时间戳不同，通过递增线性化寄存器获得的数字不会存在任何间隙。因此，如果节点完成了消息4的发送，且接收到了序列化6的消息，那么在它对消息6回复之前必须等待消息5。Lamport时间戳则不是这样，而这也是区别全序关系广播与基于时间戳排序的关键。

使用原子自增操作来创建线性化整数有多难呢？答案还是那样，如果不存在失效，就非常容易，甚至可以把它保存在某个节点的内存变量中。难点在于处理节点的网络中断，以及节点失效时如何恢复该值<sup>[59]</sup>。事实上，如果对线性化的序列号发生器深思熟虑之后所得到的最终结果，往往毫无意外地指向了共识算法。

这并非巧合，可以证明，线性化的原子比较-设置（或自增）寄存器与全序关系广播二者都等价于共识问题<sup>[28,67]</sup>。也就是说，如果你能解决其中的一个问题，那么就可以把方案用于解决其他问题。这样的结论是多么的深刻和震撼！

好了，现在终于是时候正面处理共识问题了，这是本章剩余部分的重点。

## 分布式事务与共识

共识问题是分布式计算中最重要也是最基本的问题之一。表面上看，目标只是让几个节点就某件事情达成一致。这似乎很简单，或者至少不应该太难。不幸的是，许多失败的系统正是由于低估了这个问题所导致的。





共识问题既然这么重要，本书已过大半，我们才正式揭开其面纱似乎有点姗姗来迟，这主要是因为共识主题非常之艰深，欣赏其精妙之处需要很多准备知识。即使在学术界，对共识的深刻认识也只是最近几十年才逐渐形成的，中间还纠缠着很多的误解。在探讨了复制（第5章），事务（第7章），系统模型（第8章），线性化和全序关系广播（本章）等问题之后，现在终于完成必要的准备，可以开始直面共识问题了。

有很多重要的场景都需要集群节点达成某种一致，例如：

### 主节点选举

对于主从复制的数据库，所有节点需要就谁来充当主节点达成一致。如果由于网络故障原因出现节点之间无法通信，就容易出现争议。此时，共识对于避免错误的故障切换非常重要，后者会导致两个节点都自认为是主节点即脑裂（参阅第5章“处理节点失效”）。如果集群中存在两个这样的主节点，每个都在接受写请求，最终会导致数据产生分歧、不一致甚至数据丢失。

### 原子事务提交

对于支持跨节点或跨分区事务的数据库，会面临这样的问题：某个事务可能在一些节点上执行成功，但在其他节点却不幸发生了失败。为了维护事务的原子性（即ACID，参阅第7章“原子性”），所有节点必须对事务的结果达成一致：要么全部成功提交（假定没有出错），要么中止/回滚（如果出现了错误）。这个共识的例子被称为原子提交问题<sup>注12</sup>。

## 共识的不可能性

或许你可能听说过FLP结论<sup>[68]</sup>，其字母源于三位作者Fischer, Lynch和Paterson。FLP表明如果节点存在可能崩溃的风险，则不存在总是能够达成共识的稳定算法。在分布式系统中，我们必须假设节点可能会崩溃，我们在讨论如何达成共识，而这里又说可靠的共识无法实现。到底怎么回事呢？

答案是FLP结论是基于异步系统模型而做的证明（请参阅第8章“理论系统模型与现实”），这是一个非常受限的模型，它假定确定性算法都不能使用任何时钟或超时机制。如果算法可以使用超时或其他方法来检测崩溃节点（即使怀疑可能

注12：原子提交的形式化描述其实与共识还稍有不同：原子事务只有在所有参与者都投票赞成的情况下才能完成最终提交，如果有任何参与者中止，则必须中止。共识则可以就任何一位参与者的提议进行表决，达到多数即可通过。原子提交与共识可以相互转化<sup>[70,71]</sup>。而非阻塞的原子提交则比共识更难，具体请参阅本章后面的“三阶段提交”。



是误报)，那么可以实现稳定的共识方案<sup>[67]</sup>。另外，即使算法使用了随机数来检测节点故障也可以绕过FLP结论<sup>[69]</sup>。

因此，FLP结论有其重要的理论意义，但对于实际的分布式系统通常达成共识是可行的。

本节我们将首先详细研究原子提交问题。具体来讲，我们将集中于两阶段提交（2PC）算法，这是解决原子提交最常见的方法，在各种数据库、消息系统和应用服务器中都有实现。事实证明，2PC是一种共识算法，虽然谈不上多么优秀<sup>[70,71]</sup>。

2PC学习之后，我们将继续探索更好的共识算法实现，比如ZooKeeper（Zab）和etcd（Raft）所使用的算法。

## 原子提交与两阶段提交

第7章我们了解到事务原子性的目的是，当一个包含多笔写操作的事务在执行过程出现任何意外，原子性可以为上层应用提供非常简单的语义：事务的结果要么是成功提交（所有事务的写入都是持久的），要么是中止（此时所有事务的写入都被回滚，即撤销或者丢弃）。

原子性可以防止失败的事务破坏系统，避免形成部分成功夹杂着部分失败。这对于多对象事务（参阅第7章“单对象与多对象事务操作”）和维护二级索引格外重要。每个二级索引都有与主数据不同的数据结构，因此，如果修改了某些数据，则相应的二级索引也需要随之更新。原子性可以确保二级索引与主数据总是保持一致（如果发生了不一致，那么索引的作用将会大打折扣）。

### 从单节点到分布式的原子提交

对于在单个数据库节点上执行的事务，原子性通常由存储引擎来负责。当客户端请求数据库节点提交事务时，数据库首先使事务的写入持久化（通常保存在预写日志中，请参阅第3章“可靠的B-tree”），然后把提交记录追加写入到磁盘的日志文件中。如果数据库在该过程中间发生了崩溃，那么当节点重启后，事务可以从日志中恢复；如果在崩溃之前提交记录已成功写入磁盘，则认为事务已安全提交；否则，回滚该事务的所有写入。

因此，在单节点上，事务提交非常依赖于数据持久写入磁盘的顺序关系：先写入数

据，然后再提交记录<sup>[72]</sup>。事务提交（或中止）的关键点在于磁盘完成日志记录的时刻：在完成日志记录写之前如果发生了崩溃，则事务需要中止；如果在日志写入完成之后，即使发生崩溃，事务也被安全提交。这就是在单一设备上（某个特定的磁盘连接到一个特定的节点）上实现原子提交的核心思路。

但是，如果一个事务涉及多个节点呢？例如，一个分区数据库中多对象事务，或者是基于词条分区的二级索引（其中索引条目可能位于与主数据不同的节点上，请参阅第6章“分区与二级索引”）。虽然大多数NoSQL分布式数据都不支持这样的分布式事务，但有很多集群关系型数据库则支持（参阅本章后面“实践中的分布式事务”）。

向所有节点简单地发送一个提交请求，然后各个节点独立执行事务提交是绝对不够的。这样做很容易发生部分节点提交成功，而其他一些节点发生失败，从而违反了原子性保证：

- 某些节点可能会检测到违反约束或有冲突，因而决定中止，而其他节点则可能成功提交。
- 某些提交请求可能在网络中丢失，最终由于超时而中止，而其他提交请求则顺利通过。
- 某些节点可能在日志记录写入之前发生崩溃，然后在恢复时回滚，而其他节点则成功提交。

如果一部分节点提交了事务，而其他节点却放弃了事务，节点之间就会变得不一致（见图7-3）。而且某个节点一旦提交了事务，即使事后发现其他节点发生中止，它也无法再撤销已提交的事务。正因如此，如果有部分节点提交了事务，则所有节点也必须跟着提交事务。

事务提交不可撤销，不能事后再改变主意（在提交之后再追溯去中止）。这些规则背后的深层原因是，一旦数据提交，就被其他事务可见，继而其他客户端会基于此做出相应的决策。这个原则构成了读-提交隔离级别的基础（参阅第7章“读-提交”）。如果允许事务在提交后还能中止，会违背之后所有读-提交的事务，进而被迫产生级联式的追溯和撤销。

当然已提交事务的效果可以被之后一笔新的事务来抵消掉，即补偿性事务<sup>[73,74]</sup>。不过，从数据库的角度来看，前后两个事务完全互相独立。类似这种跨事务的正确性需要由应用层来负责。



## 两阶段提交

两阶段提交（two-phase commit, 2PC）是一种在多节点之间实现事务原子提交的算法，用来确保所有节点要么全部提交，要么全部中止。它是分布式数据库中的经典算法之一<sup>[13,35,75]</sup>。2PC在某些数据库内部使用，或者以XA事务<sup>[76,77]</sup>形式（例如Java Transaction API）或SOAP Web服务WS-AtomicTransaction<sup>[78,79]</sup>的形式提供给应用程序。

2PC的基本流程如图9-9所示。不同于单节点上的请求提交，2PC中的提交/中止过程分为两个阶段（因此得名2PC）。

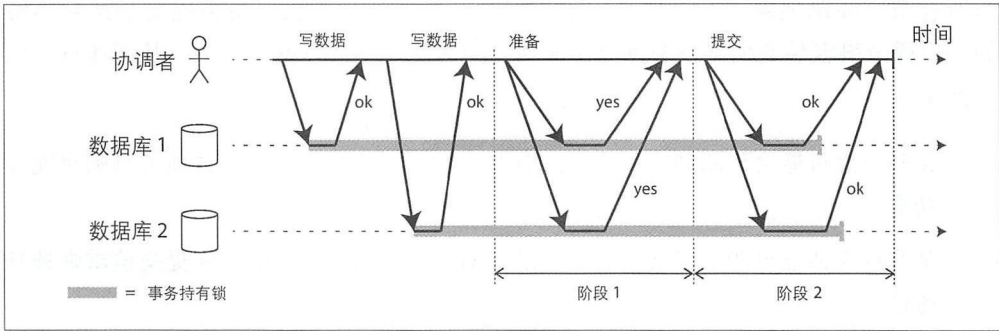


图9-9：两阶段提交（2PC）的一次成功执行



### 不要混淆2PC和2PL

两阶段提交（2PC）和两阶段加锁（参阅第7章“两阶段加锁”）是两个完全不同的事情。2PC在分布式数据库中负责原子提交，而2PL则提供可串行化的隔离。为避免混淆，最好将它们视为完全独立的概念，并忽略名称中那一点相似性。

2PC引入了单节点事务所没有的一个新组件：协调者（也称为事务管理器）。协调者通常实现为共享库，运行在请求事务相同进程中（例如嵌入在Java EE容器中），但也可以是单独的进程或服务。常见协调者的例子包括Narayana, JOTM, BTM或MSDTC。

通常，2PC事务从应用程序在多个数据库节点上执行数据读/写开始。我们将这些数据库节点称为事务中的参与者。当应用程序准备提交事务时，协调者开始阶段1：发送一个准备请求到所有节点，询问他们是否可以提交。协调者然后跟踪参与者的回应：

- 如果所有参与者回答“是”，表示他们已准备好提交，那么协调者接下来在阶段2会发出提交请求，提交开始实际执行。

- 如果有任何参与者回复“否”，则协调者在阶段2中向所有节点发送放弃请求。

这个过程有点像西方传统的婚礼仪式：主持人会询问新郎和新娘是否愿意与对方结为夫妇，通常双方都会回答“我愿意”，在确认之后和所有与会者共同见证下，方可宣布完成了婚姻承诺。而万一新郎或者新娘没有肯定回复，理论上仪式应该中止<sup>[73]</sup>。

## 系统的承诺

只从上面简单的描述可能还是不清楚为什么两阶段提交可以确保跨节点的原子性，而单一提交却做不到。试想，即使对于2PC，准备和提交请求也一样可能发生丢失，那么2PC究竟为何不同？

为了理解其工作原理，我们来更详细地分解这个过程：

1. 当应用程序启动一个分布式事务时，它首先向协调者请求事务ID。该ID全局唯一。
2. 应用程序在每个参与节点上执行单节点事务，并将全局唯一事务ID附加到事务上。此时，读写都是在单节点内完成。如果在这个阶段出现问题（例如节点崩溃或请求超时），则协调者和其他参与者都可以安全中止。
3. 当应用程序准备提交时，协调者向所有参与者发送准备请求，并附带全局事务ID。如果准备请求有任何一个发生失败或者超时，则协调者会通知所有参与者放弃事务。
4. 参与者在收到准备请求之后，确保在任何情况下都可以提交事务，包括安全地将事务数据写入磁盘（不能以任何借口稍后拒绝提交，包括系统崩溃，电源故障或磁盘空间不足等），并检查是否存在冲突或约束违规。一旦向协调者回答“是”，节点就承诺会提交事务。换句话说，尽管还没有真正提交，但参与者已表态此后不会行使放弃事务的权利。
5. 当协调者收到所有准备请求的答复时，就是否提交（或放弃）事务要做出明确的决定（即只有所有参与者都投赞成票时才会提交）。协调者把最后的决定写入到磁盘的事务日志中，防止稍后系统崩溃，并可以恢复之前的决定。这个时刻称为提交点。
6. 协调者的决定写入磁盘之后，接下来向所有参与者发送提交（或放弃）请求。如果此请求出现失败或超时，则协调者必须一直重试，直到成功为止。此时，所有节点不允许有任何反悔：开弓没有回头箭，一旦做了决定，就必须贯彻执行，即使需要很多次重试。而如果有参与者在在此期间出现故障，在其恢复之后，也必须



继续执行。这是因为之前参与者都投票选择了“是”，对于做出的承诺同样没有反悔的余地。

由此可见，该协议有两个关键的“不归路”：首先，当参与者投票“是”时，它做出了肯定提交的承诺（尽管还取决于其他的参与者的投票，协调者才能做出最后觉得）。其次，协调者做出了提交（或者放弃）的决定，这个决定也是不可撤销。正是这两个承诺确保了2PC的原子性（而单节点原子提交其实是将两个事件合二为一，写入事务日志即提交）。

回到婚姻的例子中，在说“我愿意”之前，新娘/新郎都有“放弃”承诺的自由，比如说“我不愿意！”。而在做出肯定的承诺之后，就不能随便撤销。假如在说了“我愿意”之后不巧晕倒在地，即使他/她没有听到“你们现已结为夫妻”的证词，也不会因此就能改变事实。当稍后恢复了意识，可以询问证婚人（通过事务ID状态）是否已完成婚礼，或者等待证婚人安排下一次仪式（即重试将一直持续下去）。

## 协调者发生故障

如果参与者或者网络在2PC期间发生失败，例如在第一阶段，任何一个准备请求发生了失败或者超时，那么协调者就会决定中止交易；或者在第二阶段发生提交（或中止）请求失败，则协调者将无限期重试。但是，如果协调者本身发生了故障，接下来会发生什么现在还不太清楚。

如果协调者在发送准备请求之前就已失败，则参与者可以安全地中止交易。但是，一旦参与者收到了准备请求并做了投票“是”，则参与者不能单方面放弃，它必须等待协调者的决定。如果在决定到达之前，出现协调者崩溃或网络故障，则参与者只能无奈等待。此时参与者处在一种不确定的状态。

情况如图9-10所示。在该例子中，协调者实际上做出了提交决定，数据库2已经收到了提交请求。但是，协调者在将提交请求发送到数据库1之前发生了崩溃，因此数据库1不知道该提交还是中止。超时机制也无法解决问题：如果超时之后数据库1决定单方面中止，最终将与完成提交的数据库2产生不一致。同理，参与者也不能单方面决定提交，因为可能有些参与者投了否决票导致协调者最终的决定是放弃。

没有协调者的消息，参与者无法知道下一步的行动（是提交还是放弃）。理论上，参与者之间可以互相通信，通过了解每个参与者的投票情况并最终达成一致，不过这已经不是2PC协议的范畴了。





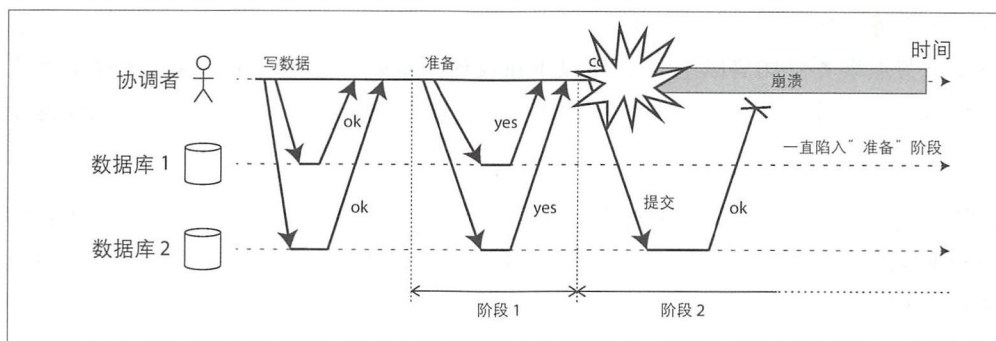


图9-10：参与者在投票成票之后，协调者发生崩溃，则数据库1不清楚是否应该提交或中止

2PC能够顺利完成的唯一方法是等待协调者恢复。这就是为什么协调者必须在向参与者发送提交（或中止）请求之前要将决定写入磁盘的事务日志：等协调者恢复之后，通过读取事务日志来确定所有未决的事务状态。如果在协调者日志中没有完成提交记录就会中止。此时，2PC的提交点现在归结为协调者在常规单节点上的原子提交。

### 三阶段提交

两阶段提交也被称为阻塞式原子提交协议，因为2PC可能在等待协调者恢复时卡住。理论上，可以使其改进为非阻塞式从而避免这种情况。但是，实践中要想做到这一点并不容易。

作为2PC的替代方案，目前也有三阶段提交算法<sup>[13, 80]</sup>。然而，3PC假定一个有界的网络延迟和节点在规定时间内响应。考虑到目前大多数具有无限网络延迟和进程暂停的实际情况（见第8章），它无法保证原子性。

通常，非阻塞原子提交依赖于一个完美的故障检测器<sup>[67,71]</sup>，即有一个非常可靠的机制可以判断出节点是否已经崩溃。在无限延迟的网络环境中，超时机制并不是可靠的故障检测器，因为即使节点正常，请求也可能由于网络问题而最终超时。正是由于这样的原因，尽管大家已经意识到上述协调者潜在的问题，但还在普遍使用2PC。

## 实践中的分布式事务

分布式事务，尤其是那些通过两阶段提交所实现的事务，声誉混杂。一方面，它们被看作是提供了一个其他方案难以企及的重要的安全保证；但另一方面，他们由于操作上的缺陷、性能问题、承诺不可靠等问题而遭受诟病<sup>[81-84]</sup>。目前，许多云服务提供商由于运维方面的问题而决定不支持分布式事务<sup>[85,86]</sup>。



分布式事务的某些实现存在严重的性能问题。例如，有报告显示MySQL的分布式事务比单节点事务慢10倍以上<sup>[87]</sup>，所以不建议使用也就不足为奇了。两阶段提交性能下降的主要原因是为了防崩溃恢复而做的磁盘I/O（fsync）<sup>[88]</sup>以及额外的网络往返开销。

但是，我们不应该就这么直接地抛弃分布式事务，而应该更加审慎的对待，从中获取一些重要的经验教训。首先，我们还是要明确“分布式事务”的确切含义。目前由两种截然不同的分布式事务概念：

### 数据库内部的分布式事务

某些分布式数据库（例如那些标配支持复制和分区的数据库）支持跨数据库节点的内部事务。例如，VoltDB和MySQL Cluster的NDB存储引擎就支持这样的内部分布式事务。此时，所有参与节点都运行着相同的数据库软件。

### 异构分布式事务

在异构分布式事务中，存在两种或两种以上不同的参与者实现技术。例如来自不同供应商的数据库，甚至是非数据库系统（如消息中间件）。即使是完全不同的系统，跨系统的分布式事务必须确保原子提交，

数据库内部事务由于不必考虑与其他系统的兼容，因此可以使用任何形式的内部协议并采取有针对性的优化。因此，数据库内部的分布式事务往往可行且工作不错，但异构环境的事务则充满了挑战。

## Exactly-once消息处理

异构的分布式事务旨在无缝集成多种不同的系统。例如，当且仅当数据库中处理消息的事务成功提交，消息队列才会标记该消息已处理完毕。这个过程是通过自动提交消息确认和数据库写入来实现的。即使消息系统和数据库两种不同的技术运行在不同的节点上，采用分布式事务也能达到上述目标。

如果消息发送或数据库事务任何一个发生失败，则两者都须中止，消息队列可以在稍后再次重传消息。因此，通过自动提交消息和消息处理的结果，可以确保消息可以有效处理有且仅有一次（成功之前有可能需要重试）。而如果事务最后发生中止，则会放弃所有部分完成的结果。

需要指出，只有在所有受影响的系统都使用相同的原子提交协议的前提下，这种分布式事务才是可行。例如，如果处理结果之一是发送一封邮件，而邮件服务器却不支持两阶段提交，此时如果某个环节出错需要重试，就会导致邮件系统重复发送两次或更多。但如果假定所有结果或者副作用都可以在事务中止时回滚，就可以安全地重新处



理消息，好像之前什么都没发生过一样。

在之后第11章我们还会回到消息的恰好一次处理问题。这里我们继续探讨异构环境下分布式事务的原子提交。

## XA交易

X/Open XA (eXtended Architecture, XA) 是异构环境下实施两阶段提交的一个工业标准<sup>[76,77]</sup>，于1991年推出并得到广泛推广。目前，许多传统关系数据库（包括PostgreSQL、MySQL、DB2、SQL Server和Oracle）和消息队列（包括ActiveMQ、HornetQ、MSMQ和IBM MQ）都支持XA。

XA并不是一个网络协议，而是一个与事务协调者进行通信的C API。当然，它也支持其他语言的API绑定，例如Java EE中，XA事务是由Java事务API（Java Transaction API, JTA）来实现，JTA可以支持非常多JDBC（Java Database Connectivity）驱动和消息队列驱动（通过Java消息服务，JMS）。

XA假定应用程序通过网络或客户端的库函数与参与者（包括数据库、消息服务）节点进行通信。如果驱动程序支持XA，意味着应用可以调用XA API来确定操作是否是异构分布式事务的一部分。如果是，则发送必要的信息给数据库服务器。它还支持回调，这样协调者可以通过回调函数来通知所有参与者执行准备或者提交（或者中止）。

事务协调者需要实现XA API。虽然标准并没有详细要求该如何实现，但实际上，协调者通常也是一个API库，它与产生事务的应用程序运行在相同的进程中。这些API会跟踪事务中所有的参与者，协调节点进行准备（通过回调）工作，然后负责收集参与者的投票，并在本地磁盘的日志文件里记录事务最终的决定。

如果应用程序进程发生崩溃，或者所在的节点出现故障，协调者就需要做相应的处理。此时，所有完成了准备阶段但尚未提交的参与者就会陷入停顿。由于事务日志保存在应用服务器的本地磁盘上，该节点必须先重启，然后协调者通过XA API读取日志、进而恢复事务的决定。完成这些之后，协调者才能继续使用数据库驱动XA回调来要求所有参与者执行提交（或中止）。数据库服务器无法直接与协调者进行通信，而须通过相应的API接口。

## 停顿时仍持有锁

为什么我们非常关注陷入停顿的参与者节点（即不确定该提交还是中止）呢？难道系统不能选择忽略（并最终清理）这些节点，这样系统不就可以继续工作么？







问题的关键在于锁。正如第7章“读-提交”所讨论的，数据库事务通常持有待修改行的行级独占锁，用以防止脏写。此外，如果要使用可串行化的隔离，则两阶段锁的数据库还会对事务曾经读取的行持有读-共享锁（参阅第7章“两阶段加锁”）。

在事务提交（或中止）之前，数据库都不会释放这些锁（图9-9中的阴影区域所示）。因此，在两阶段提交时，事务在整个停顿期间一直持有锁。换句话说，如果协调者崩溃并且需要20分钟才能重启恢复，那么这些对象将被锁定20分钟；如果协调者的日志由于某种原因而彻底丢失，这些数据对象将永久处于加锁状态，至少管理员采用手动方式解决之前只能如此。

数据处于加锁时，其他事务就无法执行修改。取决于数据库的具体实现，其他事务甚至无法读取这些行。因此，其他的事务事实上无法有效执行。这可能会导致很多上层应用基本处于不可用状态，所以必须解决处于停顿状态的那些事务。

## 从协调者故障中恢复

理论上，如果协调者崩溃之后重新启动，它应该可以从日志中恢复那些停顿的事务。然而，在实践中，孤立的不确定事务确实会发生<sup>[89,90]</sup>。无论何种原因，例如由于软件bug导致交易日志丢失或者损坏，最终协调者还是出现了恢复失败。那些悬而未决的事务无法自动解决，而是永远留在那里，而且还持有锁并阻止其他事务。

即使重启那些处于停顿状态的数据库节点也无法解决这个问题，这是由于2PC的正确实现要求即使发生了重启，也要继续保持重启之前事务的加锁（否则就会违背原子性保证）。所以，这的确非常棘手。

唯一的出路只能是让管理员手动决定究竟是执行提交还是回滚。管理员必须仔细检查每个有问题的参与者，确定是否有节点已经事实完成了提交（或中止），然后将相同的结果一一应用于所有的参与者上。这种方案可能需要大量的手工操作，而且很可能处在关键生产环境的中断间隙，背负着巨大的压力和时间限制（要不然，为什么协调者容易出现这种问题）。

许多XA的实现都支持某种紧急避险措施称之为启发式决策：这样参与者节点可以在紧急情况下单方面做出决定，放弃或者继续那些停顿的事务，而不需要等到协调者发出指令<sup>[76,77,91]</sup>。需要说明的是，这里的启发式其实是可能破坏原子性的委婉说法，它的确违背了两阶段提交所做出的承诺。因此，这种启发式决策只是为了应急，不能作为常规手段来使用。





## 分布式事务的限制

XA事务解决了多个参与者之间如何达成一致这样一个非常现实而重要的问题，但正如上面所看到的，它也引入了不少操作方面的限制。特别是，核心的事务协调者本身就是一种数据库（存储事务的投票结果），因此需要和其他重要的数据库一样格外小心：

- 如果协调者不支持数据复制，而是在单节点上运行，那么它就是整个系统的单点故障（因为它的故障导致了很多应用阻塞在停顿事务所持有的锁上）。而现实情况是，有许多协调者的实现默认情况下并非高可用，或者只支持最基本的复制。
- 许多服务器端应用程序都倾向于无状态模式（因为更受HTTP的青睐），而所有的持久状态都保存在数据库中，这样应用服务器可以轻松地添加或删除实例。但是，当协调者就是应用服务器的一部分时，部署方式就发生了根本的变化。突然间，协调者的日志成为可靠系统的重要组成部分，它要求与数据库本身一样重要（需要协调者日志恢复那些有疑问的事务）。这样的应用服务器已经不再是无状态。
- 由于XA需要与各种数据系统保持兼容，它最终其实是多系统可兼容的最低标准。例如，它无法深入检测不同系统之间的死锁条件（因为这就将需要另一个标准化协议，使得多个系统交换事务所等待的锁信息），而且不适用于SSI（参阅第7章“可串行化的快照隔离”），后者要求一个复杂的协议来识别不同系统间的写冲突。
- 对于数据库内部的分布式事务（而不是XA），限制则少很多，例如SSI的分布式版本是可行的。然而，2PC要成功提交事务还是存在潜在的限制，它要求必须所有参与者都投票赞成，如果有任何部分发生故障，整个事务只能失败。所以分布式事务有扩大事务失败的风险，这与我们构建容错系统的目标有些背道而驰。

这是否意味着我们应该放弃保持多个系统一致的希望呢？不完全是，还有其他的方法可以实现这一目标且无需担心异构分布式事务的这些负面影响，具体将在第11章和第12章呈现。现在，我们应该可以总结一下共识问题。

## 支持容错的共识

通俗理解，共识是让几个节点就某项提议达成一致。例如，多个人同时尝试预订飞机的最后一个座位或剧院中的同一座位，或者尝试使用相同的用户名注册账户，此时可以用共识算法来决定这些不相容的操作之中谁是获胜者。





共识问题通常形式化描述如下：一个或多个节点可以提议某些值，由共识算法来决定最终值。对于预约座位的例子，当多个顾客同时试图购买最后一个座位时，处理顾客请求的每个节点可以提议它所服务的顾客ID，最后的决定则是关于由哪个顾客获得座位。

在这个描述中，共识算法必须满足以下性质<sup>[25]注13</sup>：

协商一致性(Uniform agreement)

所有的节点都接受相同的决议。

诚实性(Integrity)

所有节点不能反悔，即对一项提议不能有两次决定。

合法性(Validity)

如果决定了值 $v$ ，则 $v$ 一定是由某个节点所提议的。

可终止性(Termination)

节点如果不崩溃则最终一定可以达成决议。

协商一致性和诚实性属性定义了共识的核心思想：决定一致的结果，一旦决定，就不能改变。有效性属性主要是为了排除一些无意义的方案：例如，无论什么建议，都可以有一个总是为空（NULL）的决定，虽然可以满足一致性和诚实性，但没有任何实际效果。

如果不关心容错，那么满足前三个属性很容易：可以强行指定某个节点为“独裁者”，由它做出所有的决定。但是，如果该节点失败，系统就无法继续做出任何决定。其实这就是在两阶段提交时所看到的：如果协调者失败了，那些处于不确定状态的参与者就无从知道下一步该做什么。

可终止性则引入了容错的思想。它重点强调一个共识算法不能原地空转，永远不做事，换句话说，它必须取得实质性进展。即使某些节点出现了故障，其他节点也必须最终做出决定。可终止性属于一种活性，而另外三种则属于安全性方面的属性（参阅第8章“安全性与活性”）。

上述共识的系统模型假定当某个节点发生崩溃后，节点就彻底消失，永远不再回来。可以这样设想，这不是由于软件错误，而是遭遇了地震，整个数据中心包括所有节点

---

注13：这种共识的特殊形式被称为统一共识，它相当于异步系统中，具有不可靠故障检测器的常规共识，具体请参阅文献<sup>[71]</sup>。学术文献通常采用术语进程而不是节点，我们此处用的是节点，目的是和本书其他部分保持一致。







被山体滑坡所摧毁，所以必须假设节点已经深埋于30英尺下的瓦砾之中，不可能重新上线。在这样的系统模型下，所有采取等待节点恢复的算法都无法满足终止性，特别是2PC不符合可终止性要求。

当然，如果所有的节点都崩溃了，那么无论何种算法都不可能继续做出决定。算法所能够容忍的失败次数和规模都有一定的限制。事实上，可以证明任何共识算法都需要至少大部分节点正确运行才能确保终止性<sup>[67]</sup>。而这个多数就可以安全地构成quorum（参阅第5章“读写quorum”）。

因此，可终止性的前提是，发生崩溃或者不可用的节点数必须小于半数节点。即便是多数节点出现了故障或者存在严重的网络问题，现在有很多实现的共识系统也可以满足安全属性：协商一致性，诚实性和合法性<sup>[92]</sup>。所以大规模的失效情况可能会导致系统无法处理请求，但不会破坏系统做出无效的决定。

大多数共识算法都假定系统不存在拜占庭式错误。对于“拜占庭式错误”（参阅第8章），即节点没有遵循协议（例如故意发送相互矛盾的消息），从而破坏协议的安全属性。研究表明，只要发生拜占庭故障的节点数小于三分之一<sup>[25,93]</sup>，也可以达成共识。不过篇幅所限，我们无法就此展开细节讨论。

## 共识算法与全序广播

最著名的容错式共识算法包括VSR<sup>[94,95]</sup>，Paxos<sup>[96-99]</sup>，Raft<sup>[22,100,101]</sup>和Zab<sup>[15,21,102]</sup>。这些算法存在诸多相似之处，但又不完全相同<sup>[103]</sup>。篇幅所限，本书无法详细介绍这些算法，除非你决定要自己实现一套这样的共识系统（这可能是非常不明智的做法，极具失败的可能<sup>[98,104]</sup>），否则只需了解它们共同的设计思想就足够了。

这些算法大部分其实并不是直接使用上述的形式化模型（提议并决定某个值，同时满足上面4个属性）。相反，他们是决定了一系列值，然后采用全序关系广播算法（参阅本章前面的“全序关系广播”）。

全序关系广播的要点是，消息按照相同的顺序发送到所有节点，有且只有一次。如果仔细想想，这其实相当于进行了多轮的共识过程：在每一轮，节点提出他们接下来想要发送的消息，然后决定下一个消息的全局顺序<sup>[67]</sup>。

所以，全序关系广播相当于持续的多轮共识（每一轮共识的决定对应于一条消息）：

- 由于协商一致性，所有节点决定以相同的顺序发送相同的消息。
- 由于诚实性，消息不能重复。



- 由于合法性，消息不会被破坏，也不是凭空捏造的。
- 由于可终止性，消息不会丢失。

VSR、Raft和Zab都直接采取了全序关系广播，这比重复性的一轮共识只解决一个提议更加高效。而Paxos则有对应的优化版本称之为Multi-Paxos。

## 主从复制与共识

第5章讨论了主从复制（参阅第5章“主节点与从节点”），所有的写入操作都由主节点负责，并以相同的顺序发送到从节点来保持副本更新。这不就是基本的全序关系广播么？那在主从复制时我们怎么没有考虑共识问题呢？

答案取决于如何选择主节点。如果主节点是由运营人员手动选择和配置的，那基本上就是一个独裁性质的“一致性算法”：只允许一个节点接受写入（并决定复制日志中的写入顺序），如果该节点发生故障，系统将无法写入，直到操作人员再手动配置新的节点成为主节点。这样的方案也能在实践中很好地发挥作用，但它需要人为干预才能取得进展，不满足共识的可终止性。

一些数据库支持自动选举主节点和故障切换，通过选举把某个从节点者提升为新的主节点（参阅第5章“处理节点失效”）。这样更接近容错式全序关系广播，从而达成共识。

但是，还有一个问题，我们之前曾讨论过脑裂：所有的节点都需要同意主节点，否则两个主节点会导致数据库出现不一致。因此，我们需要共识算法选出一位主节点。但是，如果这里描述的共识算法实际上是全序关系广播，且全序关系广播很像主从复制，但主从复制现在又需要选举主节点等。

看起来要选举一个新的主节点，我们首先需要有一个主节点。要解决共识，必须先处理共识。怎么摆脱这样一个奇怪的循环？

## Epoch和Quorum

目前所讨论的所有共识协议在其内部都使用了某种形式的主节点，虽然主节点并不是固定的。相反，他们都采用了一种弱化的保证：协议定义了一个世代编号（epoch number，对应于Paxos中的ballot number，VSP中view number，以及Raft中的term number），并保证在每个世代里，主节点是唯一确定的。

如果发现当前的主节点失效，节点就开始一轮投票选举新的主节点。选举会赋予一个单调递增的epoch号。如果出现了两个不同的主节点对应于不同epoch号码（例如，上



一个epoch号码的主节点其实并没有真正挂掉），则具有更高epoch号码的主节点将获胜。

在主节点做出任何决定之前，它必须首先检查是否存在比它更高的epoch号码，否则就会产生冲突的决定。主节点如何知道它是否已被其他节点所取代了呢？还记得上一章“真理由多数决定”么？节点不能依靠自己所掌握的信息来决策，例如自认为是主节点并不代表其他节点都接受了它的“自认为”。

相反，它必须从quorum节点中收集投票（参阅第5章“读写quorum”）。主节点如果想要做出某个决定，须将提议发送给其他所有节点，等待quorum节点的响应。quorum通常（但不总是）由多数节点组成<sup>[105]</sup>。并且，只有当没有发现更高epoch主节点存在时，节点才会对当前的提议（带有epoch号码）进行投票。

因此，这里面实际存在两轮不同的投票：首先是投票决定谁是主节点，然后是对主节点的提议进行投票。其中的关键一点是，参与两轮的quorum必须有重叠：如果某个提议获得通过，那么其中参与投票的节点中必须至少有一个也参加了最近一次的主节点选举<sup>[105]</sup>。换言之，如果在针对提议的投票中没有出现更高epoch号码，那么可以得出这样的结论：因为没有发生更高epoch的主节点选举，当前的主节点地位没有改变，所以可以安全地就提议进行投票。

投票过程看起来很像两阶段提交（2PC）。最大的区别是，2PC的协调者并不是依靠选举产生；另外容错共识算法只需要收到多数节点的投票结果即可通过决议，而2PC则要求每个参与者都必须做出“是”才能最终通过。此外，共识算法还定义了恢复过程，出现故障之后，通过该过程节点可以选举出新的主节点然后进入一致的状态，确保总是能够满足安全属性。所有这些差异之处都是确保共识算法正确性和容错性的关键。

## 共识的局限性

共识算法对于分布式系统来说绝对是一个巨大的突破，它为一切不确定的系统带来了明确的安全属性（一致性，完整性和有效性），此外它还可以支持容错（只要大多数节点还在工作和服务可达）。共识可以提供全序关系广播，以容错的方式实现线性化的原子操作（参阅本章前面“采用全序关系广播实现线性化存储”）。

不过，也不是所有系统都采用了共识，因为好处的背后都是有代价的。这包括：

在达成一致性决议之前，节点投票的过程是一个同步复制过程。如第5章“同步与异步复制”所述，数据库通常配置为异步复制，存在某些已提交的数据在故障切换时丢





失的风险，即使这样，很多系统还是采用异步复制（而非同步复制），原因正是为了更好的性能。

共识体系需要严格的多数节点才能运行。这意味着需要至少三个节点才能容忍一个节点发生故障（剩下的三分之二形成多数），或者需要最少五个节点来容忍两个节点故障（其余五分之三形成多数）。如果由于网络故障切断了节点之间的连接，则只有多数节点所在的分区可以继续工作，剩下的少数节点分区则处于事实上的停顿状态（参阅本章前面“线性化的代价”）。

多数共识算法假定一组固定参与投票的节点集，这意味着不能动态添加或删除节点。动态成员资格的扩展特性可以在集群中的按需调整节点数，但相比于静态的成员组成，其理解程度和接受程度要低很多。

共识系统通常依靠超时机制来检测节点失效。在网络延迟高度不确定的环境中，特别是那些跨区域分布的系统，经常由于网络延迟的原因，导致节点错误地认为主节点发生了故障。虽然这种误判并不会损害安全属性，但频繁的主节点选举显著降低了性能，系统最终会花费更多的时间和资源在选举主节点上而不是原本的服务任务。

此外，共识算法往往对网络问题特别敏感。例如，Raft已被发现存在不合理的边界条件处理<sup>[106]</sup>：如果整个网络中存在某一条网络连接持续不可靠，Raft会进入一种奇怪的状态：它不断在两个节点之间反复切换主节点，当前主节点不断被赶下台，这最终导致系统根本无法安心提供服务。其他共识算法也有类似的问题，所以面对不可靠网络，如何设计更具鲁棒性的共识算法仍然是一个开放性的研究问题。

## 成员与协调服务

ZooKeeper或etcd这样的项目通常称为“分布式键值存储”或“协调与配置服务”。从它们对外提供服务的API来看则与数据库非常相像：读取、写入对应主键的值，或者遍历主键。如果他们只是个普通数据库的话，为什么要花大力气实现一个共识算法呢？它们与其他数据库有何不同之处？

为帮助理解，我们还是先简单探讨一下通常如何使用ZooKeeper这样的服务。应用程序开发者其实很少直接使用ZooKeeper，主要因为它并非通用的数据库。绝大多数情况是通过其他很多项目来间接地依赖于Zookeeper，例如HBase，Hadoop YARN，OpenStack Nova和Kafka等都依赖于在后头运行的ZooKeeper服务。那么这些项目为什么需要它呢？

ZooKeeper和etcd主要针对保存少量、可完全载入内存的数据（虽然它们最终仍要写



入磁盘以支持持久性)而设计,所以不要用它们保存大量的数据。它们通常采用容错的全序广播算法在所有节点上复制这些数据从而实现高可靠。正如之前所讨论的,全序广播主要用来实现数据库复制:每条消息代表的是数据库写请求,然后按照相同的顺序在多个节点上应用写操作,从而达到多副本之间的一致性。

ZooKeeper的实现其实模仿了Google的Chubby分布式锁服务<sup>[14,98]</sup>,但它不仅实现了全序广播(因此实现了共识),还提供了其他很多有趣的特性。所有这些特性在构建分布式系统时格外重要:

### 线性化的原子操作

使用原子比较-设置操作,可以实现加锁服务。例如如果多个节点同时尝试执行相同的操作,则确保其中只有一个会成功。共识协议保证了操作满足原子性和线性化,即使某些节点发生故障或网络随时被中断。分布式锁通常实现为一个带有到期时间的租约,这样万一某些客户端发生故障,可以最终释放锁(参阅第8章“进程暂停”)。

### 操作全序

如第8章“主节点与锁”所述,当资源被锁或者租约保护时,需要fencing令牌来防止某些客户端由于发生进程暂停而引起锁冲突。fencing令牌确保每次加锁时数字总是单调增加。ZooKeeper在实现该功能时,采用了对所有操作执行全局排序,然后为每个操作都赋予一个单调递增的事务ID(zxid)和版本号(cversion)<sup>[15]</sup>。

### 故障检测

客户端与ZooKeeper节点维护一个长期会话,客户端会周期性地与ZooKeeper服务节点互相交换心跳信息,以检查对方是否存活。即使连接出现闪断,或者某个ZooKeeper节点发生失效,会话仍处于活动状态。但是,如果长时间心跳停止且超过了会话超时设置,ZooKeeper会声明会话失败。此时,所有该会话持有的锁资源可以配置为自动全部释放(ZooKeeper称之为ephemeral nodes即临时节点)。

### 更改通知

客户端不仅可以读取其他客户端所创建的锁和键值,还可以监视它们的变化。因此,客户端可以知道其他客户端何时加入了集群(基于它写入ZooKeeper的值)以及客户端是否发生了故障(会话超时导致节点消失)。通过订阅通知机制,客户端不需要频繁地轮询服务即可知道感兴趣对象的变化情况。

在上述特征中,其实只有线性化的原子操作才依赖于共识。然而ZooKeeper集成了所

有这些功能，在分布式协调服务中发挥了关键作用。

## 节点任务分配

ZooKeeper和Chubby系统非常适合的一个场景是，如果系统有多个流程或服务的实例，并且需求其中的一个实例充当主节点；而如果主节点失效，由其他某个节点来接管。显然，这非常吻合主从复制数据库，此外，它对于作业调度系统（或类似的有状态服务）也非常有用。

还有另一个场景，对于一些分区资源（可以是数据库，消息流，文件存储，分布式actor system等），需要决定将哪个分区分配给哪个节点。当有新节点加入集群时，需要将某些现有分区从当前节点迁移到新节点，从而实现负载动态平衡（参阅第5章“分区再平衡”）。而当节点移除或失败时，其他节点还需要接管失败节点。

上述场景中的任务，可以借助ZooKeeper中的原子操作，ephemeral nodes和通知机制来实现。如果实现无误，它可以非常方便地使应用程序达到自动故障中恢复，减少人工干预。虽然目前出现了如Apache Curator<sup>[17]</sup>等（基于ZooKeeper客户端API）提供了更高级别的封装接口，但总体上讲真正做起来并非易事。即使这样，它仍然比从头开始实现一套共识算法好很多，全新开发共识系统非常有挑战性，目前成功记录寥寥可数<sup>[107]</sup>。

应用程序最初可能只运行在单节点，之后可能最终扩展到数千节点。试图在如此庞大的集群上进行多数者投票会非常低效。ZooKeeper通常是在固定数量的节点（通常三到五个）上运行投票，可以非常高效地支持大量的客户端。因此，ZooKeeper其实提供了一种将跨节点协调服务（包括共识，操作排序和故障检测）专业外包的方式。

通常情况下，ZooKeeper所管理的数据变化非常缓慢，类似“分区7的主节点在10.1.1.23”这样的信息，其变化频率往往在分钟甚至是小时级别。它不适合保存那些应用实时运行的状态数据，后者可能每秒产生数千甚至百万次更改。如果这样，应该考虑使用其他工具（如Apache BookKeeper<sup>[108]</sup>）。

## 服务发现

此外，ZooKeeper、etcd和Consul还经常用于服务发现。例如需要某项服务时，应该连接到哪个IP地址等。在典型的云环境中，虚拟机可能会起起停停，这种动态变化的节点无法提前知道服务节点的IP地址，因此，可以这样配置服务，每当节点启动时将其网络端口信息向ZooKeeper等服务注册，然后其他人只需向ZooKeeper的注册表中询问即可。



但是，关于服务发现是否需要共识还缺乏统一认识。通过服务名称来获取IP地址传统的查询方式是基于DNS，它使用多层缓存来实现高性能与高可用性。从DNS读取肯定不满足线性化，而现实情况是，如果DNS返回的结果是过期的旧值，通常也不会产生什么大问题<sup>[109]</sup>。总体讲，DNS对于网络产生中断时服务可用性和鲁棒性更为重要一些。

即使服务发现不需要共识，但主节点选举则肯定需要。因此，如果共识系统已经明确知道哪一个是主节点，那它可以利用这些信息来帮助次级服务来发现各自的主节点。现在一些共识系统支持只读的缓存副本。这些副本异步地接收其他共识算法所达成的决议日志，但自身并不怎么参与投票，而主要是提供不需要支持线性化的读取服务。

## 会员服务

ZooKeeper等还可以看作是会员服务范畴的一部分。关于会员服务的研究历史可以追溯到20世纪80年代，它对于构建高可靠的系统（例如空中交管）非常重要<sup>[110]</sup>。

会员服务用来确定当前哪些节点处于活动状态并属于集群的有效成员。正如在第8章中所介绍的，由于无限的网络延迟，无法可靠地检测一个节点究竟是否发生了故障。但是，可以将故障检测与共识绑定在一起，让所有节点就节点的存活达成一致意见。

这里依然存在发生误判的可能性，即节点其实处于活动状态却被错误地宣判为故障。即便这样，系统就成员资格问题的决定是全体一致的，这是最重要的。例如，选举主节点的方式可能是简单地投票选择编号最小的节点，一旦节点对于当前包含哪些成员出现了不同意见，那么共识过程就无法继续。

## 小结

本章从多个不同的角度审视了一致性与共识问题。深入研究了线性化（一种流行的一致性模型）：其目标是使多副本对外看起来好像是单一副本，然后所有操作以原子方式运行，就像一个单线程程序操作变量一样。线性化的概念简单，容易理解，看起来很有吸引力，但它的主要问题在于性能，特别是在网络延迟较大的环境中。

我们接下来探讨了因果关系，因果关系对事件进行了某种排序（根据事件发生的原因-结果依赖关系）。线性化是将所有操作都放在唯一的、全局有序时间线上，而因果性则不同，它为我们提供了一个弱一致性模型：允许存在某些并发事件，所以版本历史是一个包含多个分支与合并的时间线。因果一致性避免了线性化昂贵的协调开销，且对网络延迟的敏感性要低很多。



然而，即使意识到因果顺序（例如采用Lamport时间戳），我们发现有时无法在实践中采用这种方式，在“时间戳排序还不够”一节有这样的例子：确保用户名唯一并拒绝对同一用户名的并发注册请求。如果某个节点要同意请求，则必须以某种方式查询其他节点是否存在竞争请求。这个例子最终引导我们去探究系统的共识问题。

共识意味着就某一项提议，所有节点做出一致的决定，而且决定不可撤销。通过逐一分析，事实证明，多个广泛的问题最终都可以归结为共识，并且彼此等价（这就意味着，如果找到其中一个解决方案，就可以比较容易地将其转换为其他问题的解决方案）。这些等价的问题包括：

#### 可线性化的比较-设置寄存器

寄存器需要根据当前值是否等于输入的参数，来自动决定接下来是否应该设置新值。

#### 原子事务提交

数据库需要决定是否提交或中止分布式事务。

#### 全序广播

消息系统要决定以何种顺序发送消息。

#### 锁与租约

当多个客户端争抢锁或租约时，要决定其中哪一个成功。

#### 成员/协调服务

对于失败检测器（例如超时机制），系统要决定节点的存活状态（例如基于会话超时）。

#### 唯一性约束

当多个事务在相同的主键上试图并发创建冲突资源时，约束条件要决定哪一个被允许，哪些违反约束因而必须失败。

如果系统只存在一个节点，或者愿意把所有决策功能都委托给某一个节点，那么事情就变得很简单。这和主从复制数据库的情形是一样的，即由主节点负责所有的决策事宜，正因如此，这样的数据库可以提供线性化操作、唯一性约束、完全有序的复制日志等。

然而，如果唯一的主节点发生故障，或者出现网络中断而导致主节点不可达，这样的系统就会陷入停顿状态。有以下三种基本思路来处理这种情况：

1. 系统服务停止，并等待主节点恢复。许多XA / JTA事务协调者采用了该方式。本





质上，这种方法并没有完全解决共识问题，因为它不满足终止性条件，试想如果主节点没法恢复，则系统就会永远处于停顿状态。

2. 人为介入来选择新的主节点，并重新配置系统使之生效。许多关系数据库都采用这种方法。本质上它引入了一种“上帝旨意”的共识，即在计算机系统之外由人类来决定最终命运。故障切换的速度完全取决于人类的操作，通常比计算机慢。
3. 采用算法来自动选择新的主节点。这需要一个共识算法，我们建议采用那些经过验证的共识系统来确保正确处理各种网络异常<sup>[107]</sup>。

虽然主从数据库提供了线性化操作，且在写操作粒度级别上并不依赖于共识算法，但它仍然需要共识来维护主节点角色和处理主节点变更情况。因此，某种意义上说，唯一的主节点只是其中的一步，系统在其他环节依然需要共识（虽然不那么的频繁）。好在容错算法与共识的系统可以共存，我们在本章做了简要地介绍。

ZooKeeper等工具以一种类似外包方式为应用提供了重要的共识服务、故障检测和成员服务等。虽然用起来依然有挑战，但远比自己开发共识算法要好得多（正确处理好第8章的所有问题绝非易事）。如果面临的问题最终可以归结为共识，并且还有容错需求，那么这里给的建议是采用如ZooKeeper等验证过的系统。

尽管如此，并不是每个系统都需要共识。例如无主复制和多主复制复制系统通常并不支持全局共识。正因如此，这些系统可能会发生冲突（参阅第5章“处理写冲突”），但或许也可以接受或者寻找其他方案，例如没有线性化保证时，就需要努力处理好数据多个冲突分支以及版本合并等。

本章引用了大量分布式系统理论方面的研究。虽然理论性文章和证明有时理解起来有些困难，有时甚至包含了一些不太合理的假设条件，但它们对于指导实际工作还是极具价值：例如帮助推理系统的边界在哪里；帮助发现一些违反直觉的分布式系统潜在的缺陷。如果有时间，建议仔细阅读这些参考资料。

至此我们终于完成了本书第二部分，其中包括复制（第5章），分区（第6章），事务（第7章），分布式系统失效模型（第8章）以及最后的一致性与共识（第9章）。相信我们已经奠定了坚实的理论基础，接下来第三部分我们将面向实际环境，讨论如何基于异构模块来构建强大的应用系统。

## 参考文献

- [1] Peter Bailis and Ali Ghodsi: “Eventual Consistency Today: Limitations, Extensions,







and Beyond,” *ACM Queue*, volume 11, number 3, pages 55-63, March 2013. doi:10.1145/2460276.2462076.

[2] Prince Mahajan, Lorenzo Alvisi, and Mike Dahlin: “Consistency, Availability, and Convergence,” University of Texas at Austin, Department of Computer Science, Tech Report UTCS TR-11-22, May 2011.

[3] Alex Scotti: “Adventures in Building Your Own Database,” at *All Your Base*, November 2015.

[4] Peter Bailis, Aaron Davidson, Alan Fekete, et al.: “Highly Available Transactions: Virtues and Limitations,” at *40th International Conference on Very Large Data Bases (VLDB)*, September 2014. Extended version published as pre-print arXiv:1302.0309 [cs.DB].

[5] Paolo Viotti and Marko Vukolić: “Consistency in Non-Transactional Distributed Storage Systems,” arXiv:1512.00168, 12 April 2016.

[6] Maurice P. Herlihy and Jeannette M. Wing: “Linearizability: A Correctness Condition for Concurrent Objects,” *ACM Transactions on Programming Languages and Systems (TOPLAS)*, volume 12, number 3, pages 463-492, July 1990. doi: 10.1145/78969.78972.

[7] Leslie Lamport: “On interprocess communication,” *Distributed Computing*, volume 1, number 2, pages 77-101, June 1986. doi:10.1007/BF01786228.

[8] David K. Gifford: “Information Storage in a Decentralized Computer System,” Xerox Palo Alto Research Centers, CSL-81-8, June 1981.

[9] Martin Kleppmann: “Please Stop Calling Databases CP or AP,” *martin.kleppmann.com*, May 11, 2015.

[10] Kyle Kingsbury: “Call Me Maybe: MongoDB Stale Reads,” *aphyr.com*, April 20, 2015.

[11] Kyle Kingsbury: “Computational Techniques in Knossos,” *aphyr.com*, May 17, 2014.

[12] Peter Bailis: “Linearizability Versus Serializability,” *bailis.org*, September 24, 2014.





- [13] Philip A. Bernstein, Vassos Hadzilacos, and Nathan Goodman: *Concurrency Control and Recovery in Database Systems*. Addison-Wesley, 1987. ISBN: 978-0-201-10715-9, available online at [research.microsoft.com](http://research.microsoft.com).
- [14] Mike Burrows: “The Chubby Lock Service for Loosely-Coupled Distributed Systems,” at *7th USENIX Symposium on Operating System Design and Implementation* (OSDI), November 2006.
- [15] Flavio P. Junqueira and Benjamin Reed: *ZooKeeper: Distributed Process Coordination*. O’Reilly Media, 2013. ISBN: 978-1-449-36130-3.
- [16] “etcd 2.0.12 Documentation,” CoreOS, Inc., 2015.
- [17] “Apache Curator,” Apache Software Foundation, [curator.apache.org](http://curator.apache.org), 2015.
- [18] Morali Vallath: *Oracle 10g RAC Grid, Services & Clustering*. Elsevier Digital Press, 2006. ISBN: 978-1-555-58321-7.
- [19] Peter Bailis, Alan Fekete, Michael J Franklin, et al.: “Coordination-Avoiding Database Systems,” *Proceedings of the VLDB Endowment*, volume 8, number 3, pages 185-196, November 2014.
- [20] Kyle Kingsbury: “Call Me Maybe: etcd and Consul,” [aphyr.com](http://aphyr.com), June 9, 2014.
- [21] Flavio P. Junqueira, Benjamin C. Reed, and Marco Serafini: “Zab: High-Performance Broadcast for Primary-Backup Systems,” at *41st IEEE International Conference on Dependable Systems and Networks (DSN)*, June 2011. doi:10.1109/DSN.2011.5958223.
- [22] Diego Ongaro and John K. Ousterhout: “In Search of an Understandable Consensus Algorithm (Extended Version),” at *USENIX Annual Technical Conference (ATC)*, June 2014.
- [23] Hagit Attiya, Amotz Bar-Noy, and Danny Dolev: “Sharing Memory Robustly in Message-Passing Systems,” *Journal of the ACM*, volume 42, number 1, pages 124-142, January 1995. doi:10.1145/200836.200869.
- [24] Nancy Lynch and Alex Shvartsman: “Robust Emulation of Shared Memory Using





Dynamic Quorum-Acknowledged Broadcasts,” at *27th Annual International Symposium on Fault-Tolerant Computing (FTCS)*, June 1997. doi:10.1109/FTCS. 1997.614100.

[25] Christian Cachin, Rachid Guerraoui, and Luís Rodrigues: *Introduction to Reliable and Secure Distributed Programming*, 2nd edition. Springer, 2011. ISBN: 978-3-642-15259-7, doi:10.1007/978-3-642-15260-3.

[26] Sam Elliott, Mark Allen, and Martin Kleppmann: personal communication, thread on *twitter.com*, October 15, 2015.

[27] Niklas Ekström, Mikhail Panchenko, and Jonathan Ellis: “Possible Issue with Read Repair?,” email thread on *cassandra-dev* mailing list, October 2012.

[28] Maurice P. Herlihy: “Wait-Free Synchronization,” *ACM Transactions on Programming Languages and Systems (TOPLAS)*, volume 13, number 1, pages 124-149, January 1991. doi:10.1145/114005.102808.

[29] Armando Fox and Eric A. Brewer: “Harvest, Yield, and Scalable Tolerant Systems,” at *7th Workshop on Hot Topics in Operating Systems (HotOS)*, March 1999. doi:10.1109/HOTOS.1999.798396.

[30] Seth Gilbert and Nancy Lynch: “Brewer’s Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services,” *ACM SIGACT News*, volume 33, number 2, pages 51-59, June 2002. doi:10.1145/564585.564601.

[31] Seth Gilbert and Nancy Lynch: “Perspectives on the CAP Theorem,” *IEEE Computer Magazine*, volume 45, number 2, pages 30-36, February 2012. doi:10.1109/MC.2011.389.

[32] Eric A. Brewer: “CAP Twelve Years Later: How the ‘Rules’ Have Changed,” *IEEE Computer Magazine*, volume 45, number 2, pages 23-29, February 2012. doi:10.1109/MC.2012.37.

[33] Susan B. Davidson, Hector Garcia-Molina, and Dale Skeen: “Consistency in Partitioned Networks,” *ACM Computing Surveys*, volume 17, number 3, pages 341-370, September 1985. doi:10.1145/5505.5508.

[34] Paul R. Johnson and Robert H. Thomas: “RFC 677: The Maintenance of Duplicate Databases,” Network Working Group, January 27, 1975.







- [35] Bruce G. Lindsay, Patricia Griffiths Selinger, C. Galtieri, et al.: “Notes on Distributed Databases,” IBM Research, Research Report RJ2571(33471), July 1979.
- [36] Michael J. Fischer and Alan Michael: “Sacrificing Serializability to Attain High Availability of Data in an Unreliable Network,” at *1st ACM Symposium on Principles of Database Systems* (PODS), March 1982. doi:10.1145/588111.588124.
- [37] Eric A. Brewer: “NoSQL: Past, Present, Future,” at *QCon San Francisco*, November 2012.
- [38] Henry Robinson: “CAP Confusion: Problems with ‘Partition Tolerance,’ ” *blog. cloudera.com*, April 26, 2010.
- [39] Adrian Cockcroft: “Migrating to Microservices,” at *QCon London*, March 2014.
- [40] Martin Kleppmann: “A Critique of the CAP Theorem,” arXiv:1509.05393, September 17, 2015.
- [41] Nancy A. Lynch: “A Hundred Impossibility Proofs for Distributed Computing,” at *8th ACM Symposium on Principles of Distributed Computing* (PODC), August 1989. doi:10.1145/72981.72982.
- [42] Hagit Attiya, Faith Ellen, and Adam Morrison: “Limitations of Highly-Available Eventually-Consistent Data Stores,” at *ACM Symposium on Principles of Distributed Computing* (PODC), July 2015. doi:10.1145/2767386.2767419.
- [43] Peter Sewell, Susmit Sarkar, Scott Owens, et al.: “x86-TSO: A Rigorous and Usable Programmer’s Model for x86 Multiprocessors,” *Communications of the ACM*, volume 53, number 7, pages 89-97, July 2010. doi:10.1145/1785414.1785443.
- [44] Martin Thompson: “Memory Barriers/Fences,” *mechanicalsympathy.blogspot.co.uk*, July 24, 2011.
- [45] Ulrich Drepper: “What Every Programmer Should Know About Memory,” *akkadia.org*, November 21, 2007.
- [46] Daniel J. Abadi: “Consistency Tradeoffs in Modern Distributed Database System Design,” *IEEE Computer Magazine*, volume 45, number 2, pages 37-42, February 2012. doi:10.1109/MC.2012.33.





- [47] Hagit Attiya and Jennifer L. Welch: “Sequential Consistency Versus Linearizability,” *ACM Transactions on Computer Systems (TOCS)*, volume 12, number 2, pages 91-122, May 1994. doi:10.1145/176575.176576.
- [48] Mustaque Ahamad, Gil Neiger, James E. Burns, et al.: “Causal Memory: Definitions, Implementation, and Programming,” *Distributed Computing*, volume 9, number 1, pages 37-49, March 1995. doi:10.1007/BF01784241.
- [49] Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, and David G. Andersen: “Stronger Semantics for Low-Latency Geo-Replicated Storage,” at *10th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, April 2013.
- [50] Marek Zawirski, Annette Bieniusa, Valter Balegas, et al.: “SwiftCloud: Fault-Tolerant Geo-Replication Integrated All the Way to the Client Machine,” INRIA Research Report 8347, August 2013.
- [51] Peter Bailis, Ali Ghodsi, Joseph M Hellerstein, and Ion Stoica: “Bolt-on Causal Consistency,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2013.
- [52] Philippe Ajoux, Nathan Bronson, Sanjeev Kumar, et al.: “Challenges to Adopting Stronger Consistency at Scale,” at *15th USENIX Workshop on Hot Topics in Operating Systems (HotOS)*, May 2015.
- [53] Peter Bailis: “Causality Is Expensive (and What to Do About It),” *bailis.org*, February 5, 2014.
- [54] Ricardo Gonçalves, Paulo Sérgio Almeida, Carlos Baquero, and Victor Fonte: “Concise Server-Wide Causality Management for Eventually Consistent Data Stores,” at *15th IFIP International Conference on Distributed Applications and Interoperable Systems (DAIS)*, June 2015. doi:10.1007/978-3-319-19129-4\_6.
- [55] Rob Conery: “A Better ID Generator for PostgreSQL,” *rob.conery.io*, May 29, 2014.
- [56] Leslie Lamport: “Time, Clocks, and the Ordering of Events in a Distributed System,” *Communications of the ACM*, volume 21, number 7, pages 558-565, July 1978. doi:10.1145/359545.359563.





[57] Xavier Défago, André Schiper, and Péter Urbán: “Total Order Broadcast and Multicast Algorithms: Taxonomy and Survey,” *ACM Computing Surveys*, volume 36, number 4, pages 372-421, December 2004. doi:10.1145/1041680.1041682.

[58] Hagit Attiya and Jennifer Welch: *Distributed Computing: Fundamentals, Simulations and Advanced Topics*, 2nd edition. John Wiley & Sons, 2004. ISBN: 978-0-471-45324-6, doi:10.1002/0471478210.

[59] Mahesh Balakrishnan, Dahlia Malkhi, Vijayan Prabhakaran, et al.: “CORFU: A Shared Log Design for Flash Clusters,” at *9th USENIX Symposium on Networked Systems Design and Implementation* (NSDI), April 2012.

[60] Fred B. Schneider: “Implementing Fault-Tolerant Services Using the State Machine Approach: A Tutorial,” *ACM Computing Surveys*, volume 22, number 4, pages 299-319, December 1990.

[61] Alexander Thomson, Thaddeus Diamond, Shu-Chun Weng, et al.: “Calvin: Fast Distributed Transactions for Partitioned Database Systems,” at *ACM International Conference on Management of Data* (SIGMOD), May 2012.

[62] Mahesh Balakrishnan, Dahlia Malkhi, Ted Wobber, et al.: “Tango: Distributed Data Structures over a Shared Log,” at *24th ACM Symposium on Operating Systems Principles* (SOSP), November 2013. doi:10.1145/2517349.2522732.

[63] Robbert van Renesse and Fred B. Schneider: “Chain Replication for Supporting High Throughput and Availability,” at *6th USENIX Symposium on Operating System Design and Implementation* (OSDI), December 2004.

[64] Leslie Lamport: “How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs,” *IEEE Transactions on Computers*, volume 28, number 9, pages 690-691, September 1979. doi:10.1109/TC.1979.1675439.

[65] Enis Söztutar, Devaraj Das, and Carter Shanklin: “Apache HBase High Availability at the Next Level,” *hortonworks.com*, January 22, 2015.

[66] Brian F Cooper, Raghu Ramakrishnan, Utkarsh Srivastava, et al.: “PNUTS: Yahoo!’s Hosted Data Serving Platform,” at *34th International Conference on Very Large Data Bases* (VLDB), August 2008. doi:10.14778/1454159.1454167.



- [67] Tushar Deepak Chandra and Sam Toueg: “Unreliable Failure Detectors for Reliable Distributed Systems,” *Journal of the ACM*, volume 43, number 2, pages 225-267, March 1996. doi:10.1145/226643.226647.
- [68] Michael J. Fischer, Nancy Lynch, and Michael S. Paterson: “Impossibility of Distributed Consensus with One Faulty Process,” *Journal of the ACM*, volume 32, number 2, pages 374-382, April 1985. doi:10.1145/3149.214121.
- [69] Michael Ben-Or: “Another Advantage of Free Choice: Completely Asynchronous Agreement Protocols,” at *2nd ACM Symposium on Principles of Distributed Computing (PODC)*, August 1983. doi:10.1145/800221.806707.
- [70] Jim N. Gray and Leslie Lamport: “Consensus on Transaction Commit,” *ACM Transactions on Database Systems (TODS)*, volume 31, number 1, pages 133-160, March 2006. doi:10.1145/1132863.1132867.
- [71] Rachid Guerraoui: “Revisiting the Relationship Between Non-Blocking Atomic Commitment and Consensus,” at *9th International Workshop on Distributed Algorithms (WDAG)*, September 1995. doi:10.1007/BFb0022140.
- [72] Thanumalayan Sankaranarayanan Pillai, Vijay Chidambaram, Ramnatthan Alagappan, et al.: “All File Systems Are Not Created Equal: On the Complexity of Crafting Crash-Consistent Applications,” at *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, October 2014.
- [73] Jim Gray: “The Transaction Concept: Virtues and Limitations,” at *7th International Conference on Very Large Data Bases (VLDB)*, September 1981.
- [74] Hector Garcia-Molina and Kenneth Salem: “Sagas,” at *ACM International Conference on Management of Data (SIGMOD)*, May 1987. doi:10.1145/38713.38742.
- [75] C. Mohan, Bruce G. Lindsay, and Ron Obermarck: “Transaction Management in the R\* Distributed Database Management System,” *ACM Transactions on Database Systems*, volume 11, number 4, pages 378-396, December 1986. doi: 10.1145/7239.7266.
- [76] “Distributed Transaction Processing: The XA Specification,” X/Open Company Ltd., Technical Standard XO/CAE/91/300, December 1991. ISBN: 978-1-872-63024-3.
- [77] Mike Spille: “XA Exposed, Part II,” *jroller.com*, April 3, 2004.

- [78] Ivan Silva Neto and Francisco Reverbel: “Lessons Learned from Implementing WS-Coordination and WS-AtomicTransaction,” at *7th IEEE/ACIS International Conference on Computer and Information Science (ICIS)*, May 2008. doi:10.1109/ICIS.2008.75.
- [79] James E. Johnson, David E. Langworthy, Leslie Lamport, and Friedrich H. Vogt: “Formal Specification of a Web Services Protocol,” at *1st International Workshop on Web Services and Formal Methods (WS-FM)*, February 2004. doi:10.1016/j.entcs.2004.02.022.
- [80] Dale Skeen: “Nonblocking Commit Protocols,” at *ACM International Conference on Management of Data (SIGMOD)*, April 1981. doi:10.1145/582318.582339.
- [81] Gregor Hohpe: “Your Coffee Shop Doesn’t Use Two-Phase Commit,” *IEEE Software*, volume 22, number 2, pages 64-66, March 2005. doi:10.1109/MS.2005.52.
- [82] Pat Helland: “Life Beyond Distributed Transactions: An Apostate’s Opinion,” at *3rd Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2007.
- [83] Jonathan Oliver: “My Beef with MSDTC and Two-Phase Commits,” *blog.jonathanoliver.com*, April 4, 2011.
- [84] Oren Eini (Ahende Rahien): “The Fallacy of Distributed Transactions,” *ayende.com*, July 17, 2014.
- [85] Clemens Vasters: “Transactions in Windows Azure (with Service Bus) -An Email Discussion,” *vasters.com*, July 30, 2012.
- [86] “Understanding Transactionality in Azure,” *NServiceBus Documentation*, Particular Software, 2015.
- [87] Randy Wigginton, Ryan Lowe, Marcos Albe, and Fernando Ipar: “Distributed Transactions in MySQL,” at *MySQL Conference and Expo*, April 2013.
- [88] Mike Spille: “XA Exposed, Part I,” *jroller.com*, April 3, 2004.
- [89] Ajmer Dhariwal: “Orphaned MSDTC Transactions (-2 spids),” *eraofdata.com*, December 12, 2008.
- [90] Paul Randal: “Real World Story of DBCC PAGE Saving the Day,” *sqlskills.com*, June 19, 2013.

- [91] “in-doubt xact resolution Server Configuration Option,” SQL Server 2016 documentation, Microsoft, Inc., 2016.
- [92] Cynthia Dwork, Nancy Lynch, and Larry Stockmeyer: “Consensus in the Presence of Partial Synchrony,” *Journal of the ACM*, volume 35, number 2, pages 288-323, April 1988. doi:10.1145/42282.42283.
- [93] Miguel Castro and Barbara H. Liskov: “Practical Byzantine Fault Tolerance and Proactive Recovery,” *ACM Transactions on Computer Systems*, volume 20, number 4, pages 396-461, November 2002. doi:10.1145/571637.571640.
- [94] Brian M. Oki and Barbara H. Liskov: “Viewstamped Replication: A New Primary Copy Method to Support Highly-Available Distributed Systems,” at *7th ACM Symposium on Principles of Distributed Computing* (PODC), August 1988. doi: 10.1145/62546.62549.
- [95] Barbara H. Liskov and James Cowling: “Viewstamped Replication Revisited” , Massachusetts Institute of Technology, Tech Report MIT-CSAIL-TR-2012-021, July 2012.
- [96] Leslie Lamport: “The Part-Time Parliament,” *ACM Transactions on Computer Systems*, volume 16, number 2, pages 133-169, May 1998. doi:10.1145/279227.279229.
- [97] Leslie Lamport: “Paxos Made Simple,” *ACM SIGACT News*, volume 32, number 4, pages 51-58, December 2001.
- [98] Tushar Deepak Chandra, Robert Griesemer, and Joshua Redstone: “Paxos Made Live-An Engineering Perspective,” at *26th ACM Symposium on Principles of Distributed Computing* (PODC), June 2007.
- [99] Robbert van Renesse: “Paxos Made Moderately Complex,” *cs.cornell.edu*, March 2011.
- [100] Diego Ongaro: “Consensus: Bridging Theory and Practice,” PhD Thesis, Stanford University, August 2014.
- [101] Heidi Howard, Malte Schwarzkopf, Anil Madhavapeddy, and Jon Crowcroft: “Raft Refloated: Do We Have Consensus?,” *ACM SIGOPS Operating Systems Review*, volume 49, number 1, pages 12-21, January 2015. doi:10.1145/2723872.2723876.
- [102] André Medeiros: “ZooKeeper’s Atomic Broadcast Protocol: Theory and



Practice,” Aalto University School of Science, March 20, 2012.

[103] Robbert van Renesse, Nicolas Schiper, and Fred B. Schneider: “Vive La Différence: Paxos vs. Viewstamped Replication vs. Zab,” *IEEE Transactions on Dependable and Secure Computing*, volume 12, number 4, pages 472-484, September 2014. doi: 10.1109/TDSC.2014.2355848.

[104] Will Portnoy: “Lessons Learned from Implementing Paxos,” *blog.willportnoy.com*, June 14, 2012.

[105] Heidi Howard, Dahlia Malkhi, and Alexander Spiegelman: “Flexible Paxos: Quorum Intersection Revisited,” *arXiv:1608.06696*, August 24, 2016.

[106] Heidi Howard and Jon Crowcroft: “Coracle: Evaluating Consensus at the Internet Edge,” at *Annual Conference of the ACM Special Interest Group on Data Communication (SIGCOMM)*, August 2015. doi:10.1145/2829988.2790010.

[107] Kyle Kingsbury: “Call Me Maybe: Elasticsearch 1.5.0,” *aphyr.com*, April 27, 2015.

[108] Ivan Kelly: “BookKeeper Tutorial,” *github.com*, October 2014.

[109] Camille Fournier: “Consensus Systems for the Skeptical Architect,” at *Craft Conference*, Budapest, Hungary, April 2015.

[110] Kenneth P. Birman: “A History of the Virtual Synchrony Replication Model,” in *Replication: Theory and Practice*, Springer LNCS volume 5959, chapter 6, pages 91-120, 2010. ISBN: 978-3-642-11293-5, doi:10.1007/978-3-642-11294-2\_6.

# 派生数据

在本书的第一部分和第二部分中，我们就分布式数据库的所有主要注意事项进行了汇总，包括磁盘数据布局，出现故障时的分布式一致性等。但是，这些讨论都假定应用程序只包含一个数据库。

事实上，数据系统往往要更复杂。在大型应用程序中，通常需要以多种方式访问和处理数据，并且一个数据库往往无法同时满足所有不同的需求。因此，应用程序需要使用若干个不同的数据存储区、索引、高速缓存、分析系统等的组合，并实现相应机制将数据从一个存储系统移动到另一个存储系统。

在本书的最后一部分，我们将讨论如何将多个不同数据系统（可能具有不同数据模型，并针对不同访问模式进行过优化）整合至一致的应用程序体系结构中。很多厂商声称自己的产品能够满足所有的需求，却往往忽略了这一方面的系统建设。实际上，整合不同系统是大型应用中最为关键的任务之一。

## 记录系统与派生数据系统

存储与处理数据的系统按照高层次分类可以分为两大类：

### 记录系统

一个记录系统也被称为真实数据系统，拥有数据的权威版本。当新数据进入时，例如用户输入，首先被写入记录系统。每个记录在系统中只表示一次（通常为规范化表示）。如果另一个系统与记录系统之间存在任何差异，那么以记录系统中的数据值（按照定义）为准。

## 派生数据系统

派生数据系统中的数据则是从另一个系统中获取已有数据并以某种方式进行转换或处理的结果。如果派生数据丢失，用户可以从原始数据源进行重建。一个典型的例子是缓存：如果数据存在于缓存中，则可以从缓存系统读取。但是如果缓存中并不包含所需要的数据，可以回溯到底层数据库。非规范化数值、索引和物化视图也都属于派生数据类别。而在推荐系统中，预测性的汇总数据通常源于使用日志。

从技术上说，派生数据是冗余的，也就是说它是对现有信息的复制。但是，派生数据对于获得良好的读取查询性能通常至关重要。它通常是非规范化数据。可以从单一数据源派生出多个不同的数据集，使得用户能够从不同的“视角”来审视数据。

并非所有系统都会在架构中明确区分记录系统和派生数据系统。但做出这样的区分往往非常有帮助，由此来确定系统中的数据流：例如明确系统中的某部分具有哪些输入和哪些输出，以及它们之间如何相互依赖。

大多数数据库、存储引擎和查询语言本身并不属于记录系统或派生系统。数据库只是一个工具：如何使用完全取决于用户。即记录系统和派生数据系统之间的区别并不在于工具本身，而是在于如何在应用程序中使用它们。

通过弄清楚数据的来龙去脉，来帮助厘清复杂的系统架构。这一点将贯穿本书的这一部分。

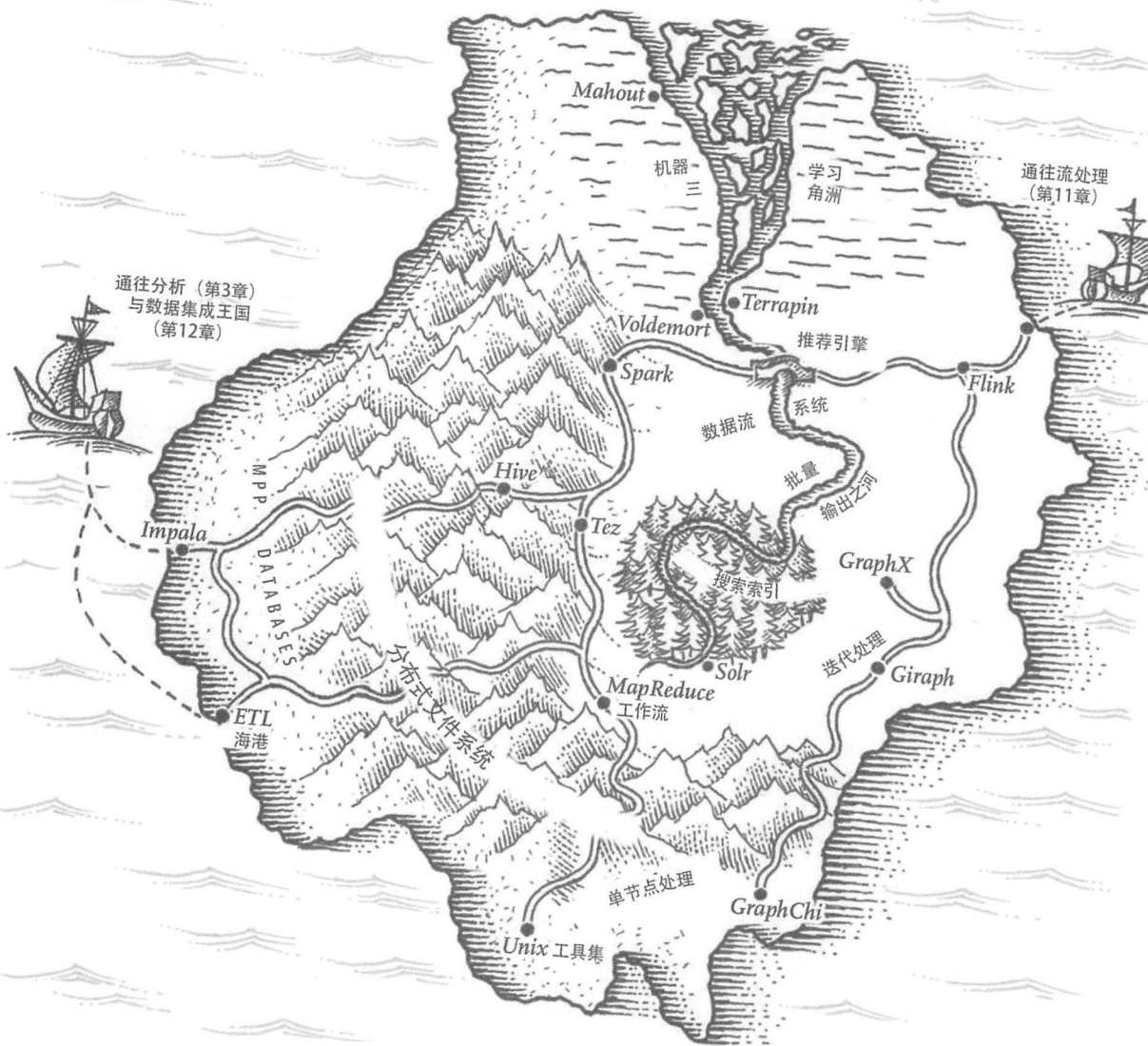
## 章节概述

我们将在第10章中讨论面向批处理的数据流系统（例如MapReduce）组件和原理，来看看它们如何提供构建大规模数据系统。第11章我们将把这些想法应用到流数据中，从而以更低的延迟完成同样的事情。第12章对本书进行了总结，探讨了未来如何使用这些工具来构建可靠、可扩展与可维护的应用程序。





派生数据之海





# 批处理系统

一个系统如果受到个人影响太大，这个系统就不可能成功。一旦最初的设计完成并且足够健壮，那么真正的测试就开始于许多持不同观点的人进行他们各自的实验。

——Donald Knuth

本书前两部分讨论了很多有关请求和查询以及对应的响应或结果等方面的内容。许多现代数据系统都假设以这种方式来处理数据：用户请求某种信息，或者发送指令，一段时间后（希望）系统会返回结果。数据库、高速缓存、搜索索引、Web服务器和其他许多系统都是这样工作的。

对于这种在线系统，无论是请求页面的浏览器还是调用远程API的服务，往往都假定请求是由用户触发，而且用户正在等待响应。通常等待不应太久，所以我们非常重视这些系统的响应时间（参阅第1章“描述性能”）。

Web和越来越多基于HTTP/REST的API使得请求 / 响应的交互模式变得如此普遍，以至于很容易将其视为理所当然。但是我们应当记住，这并不是构建系统的唯一途径，其他方法也有其优点。下面我们来区分三种不同类型的系统：

### 在线服务（或称在线系统）

服务等待客户请求或指令的到达。当收到请求或指令时，服务试图尽可能快地处理它，并发回一个响应。响应时间通常是服务性能的主要衡量指标，而可用性同样非常重要（如果客户端无法访问服务，用户可能会收到一个报错消息）。

### 批处理系统（或称离线系统）

批处理系统接收大量的输入数据，运行一个作业来处理数据，并产生输出数据。





作业往往需要执行一段时间（从几分钟到几天），所以用户通常不会等待作业完成。相反，批量作业通常会定期运行（例如，每天一次）。批处理作业的主要性能衡量标准通常是吞吐量（处理一定大小的输入数据集所需的时间）。本章主要讨论批处理。

### 流处理系统（或称近实时系统）

流处理介于在线与离线/批处理之间（所以有时称为近实时或近线处理）。与批处理系统类似，流处理系统处理输入并产生输出（而不是响应请求）。但是，流式作业在事件发生后不久即可对事件进行处理，而批处理作业则使用固定的一组输入数据进行操作。这种差异使得流处理系统比批处理系统具有更低的延迟。由于流处理是在批处理的基础上进行的，我们将在第11章讨论流处理系统。

正如我们将在本章中所看到的，批处理是构建可靠、可扩展与可维护应用的重要组成部分。例如，2004年发表的著名批处理算法MapReduce<sup>[1]</sup>“使得Google具有如此大规模的可扩展性能力”（虽然该说法有些夸大和简单化）<sup>[2]</sup>。该算法随后在各种开源数据系统中被陆续实现，包括Hadoop、CouchDB和MongoDB。

与多年前为数据仓库开发的并行处理系统相比，MapReduce是一个相当低级别的编程模型<sup>[3, 4]</sup>，但是对于运行在商用硬件上的处理规模来讲，它绝对是一个重大的进步。虽然MapReduce的重要性现在有些下降<sup>[5]</sup>，但它仍然值得深入理解，因为它清晰地解释了批处理为什么有用以及如何有用。

事实上，批处理是一种非常古老的计算形式。早在可编程数字计算机诞生之前，打孔卡制表机（例如1890年美国人口普查<sup>[6]</sup>中使用的Hollerith机器）就实现了半机械化的批量处理，以计算来自大量输入的汇总统计信息。而MapReduce与1940年和1950年广泛应用于商业数据处理的IBM卡片分类机器有着惊人的相似之处<sup>[7]</sup>。是的，历史总是在重演。

本章将介绍MapReduce和其他一些批处理算法和框架，并探讨它们在现代数据系统中的应用。但首先，我们从使用标准UNIX工具的数据处理开始。或许你已经很熟悉，但考虑到UNIX的思想和经验普遍见于大规模、异构分布式数据系统，因此值得重温UNIX哲学。

## 使用UNIX工具进行批处理

我们从一个简单的例子开始。假设有一个Web服务器，每次响应请求时都会在日志文件中追加一行记录。例如，使用nginx默认访问日志格式，日志中的一行记录如下所示：







```
216.58.210.78 - - [27/Feb/2015:17:55:11 +0000] "GET /css/typography.css HTTP/1.1"
200 3377 "http://martin.kleppmann.com/" "Mozilla/5.0 (Macintosh; Intel Mac OS X
10_9_5) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/40.0.2214.115
Safari/537.36"
```

实际上这只是一行，分成多行只是为了方便阅读。这一行中包含很多信息。为了解释它，你需要了解日志格式的定义，如下所示：

```
$remote_addr - $remote_user [$time_local] "$request"
$status $body_bytes_sent "$http_referer" "$http_user_agent"
```

因此，这一行日志表示2015年2月27日17:55:11 UTC，服务器从IP地址为216.58.210.78的客户端接收到一个对文件/css/typography.css的请求。用户为非认证用户，所以\$remote\_addr被设置为连字符（-）。响应状态是200（即请求成功），响应大小是3377字节。浏览器是Chrome 40，并且它加载了该文件，因为这个文件在URL为http://martin.kleppmann.com/的页面中被引用。

## 简单日志分析

有很多工具可以用来处理日志文件，并生成关于网站流量的漂亮报告。但是出于实践的目的，我们从基本的UNIX工具做起。例如，假设有找出网站中前五个最受欢迎的网页，可以在UNIX shell中执行下列操作<sup>注1</sup>：

```
cat /var/log/nginx/access.log | ❶
awk '{print $7}' | ❷
sort | ❸
uniq -c | ❹
sort -r -n | ❺
head -n 5 ❻
```

- ❶ 读取日志文件。
- ❷ 将每一行按空格分割成不同的字段，每行只输出第七个字段，即请求的URL地址。在我们的示例行中，这个请求URL地址是/css/typography.css。
- ❸ 按字母顺序排列URL地址列表。如果某个URL被请求过*n*次，那么排序后，结果中将包含连续*n*次的重复URL。
- ❹ uniq命令通过检查两条相邻的行是否相同来过滤掉其输入中的重复行。-c选项为输出一个计数器：对于每个不同的URL，它会报告输入中出现该URL的次数。

---

注1：有些人可能会指出cat在这里是不必要的，因为输入文件可以直接作为awk的参数给出。不过，这种写法的线性流水线特征更为明显。





- ⑤ 第二种排序按每行起始处的数字（-n）排序，也就是URL的请求次数。然后以反向（-r）顺序输出结果，即结果中最大的数字首先返回。
- ⑥ 最后，head只输出输入中的前五（-n 5），并丢弃其他数据。

该命令序列的输出如下所示：

```
4189 /favicon.ico
3631 /2013/05/24/improving-security-of-ssh-private-keys.html
2124 /2012/12/05/schema-evolution-in-avro-protocol-buffers-thrift.html
1369 /
915 /css/typography.css
```

如果你不熟悉UNIX工具，上面的命令行可能看起来令人费解，但是它的确非常强大。它将在几秒钟内处理千兆字节的日志文件，你可以轻松修改分析指令以满足自己的需求。例如，如果要省略CSS文件，将awk参数更改为'\$7 !~ /\.css\$/ {print \$7}'即可。如果想得到请求次数最多的客户端IP地址而不是页面，那么就将awk参数更改为'{print \$1}'。

本书并不打算详细探讨UNIX工具，但是这些工具非常值得学习。令人惊讶的是，使用awk, sed, grep, sort, uniq和xargs的组合，可以在几分钟内完成许多数据分析任务，并且表现令人十分满意<sup>[8]</sup>。

## 命令链与自定义程序

你也可以写一个简单的程序来做同样的事情，而不是使用UNIX命令链。例如，Ruby代码的实现可能看起来像这样：

```
counts = Hash.new(0) ❶

File.open('/var/log/nginx/access.log') do |file|
  file.each do |line|
    url = line.split[6] ❷
    counts[url] += 1 ❸
  end
end

top5 = counts.map{|url, count| [count, url] }.sort.reverse[0...5] ❹
top5.each{|count, url| puts "#{count} #{url}" } ❺
```

- ❶ counts是一个哈希表，保存了用于记录每个URL出现次数的计数器。计数器默认值为0。
- ❷ 对日志的每一行，都把URL作为第七个被空格分隔的字段（这里的数组索引是6，因为Ruby数组是零索引的）。





- ③ 为日志当前行中URL地址的计数器加1。
- ④ 按计数器值（降序）对哈希表内容进行排序，并取前五位。
- ⑤ 打印前五个条目。

这个程序并不像UNIX管道那样简洁，但是它的可读性很强，所以喜欢哪一个完全是个人爱好。除了表面上的语法差异之外，执行流程也存在很大差异。如果在大文件上运行此分析，差异会变得更加明显。

## 排序与内存中聚合

Ruby脚本需要一个URL的内存哈希表，其中每个URL地址都会映射到被访问的次数。UNIX流水线例子中则没有这样的哈希表，而是依赖于对URL列表进行排序，在这个URL列表中多次出现的相同URL仅仅是简单重复。

哪种方法更好呢？这取决于有多少个不同的网址。对于大多数中小型网站，也许可以在内存中（比如说1GB内存）存储所有不同的URL，并且可以为每个URL提供一个计数器。在此示例中，作业的工作集（作业需要随机访问的内存量）仅取决于不同URL的数量：如果单个URL有一百万条日志条目，则哈希表中所需的空间表仍然只是一个URL加上计数器的大小。如果这个工作集足够小，那么内存哈希表工作正常，甚至在笔记本计算机上都能工作。

另一方面，如果作业的工作集大于可用内存，则排序方法的优点是可以高效地使用磁盘。这与第3章的“SSTables和LSM-Tree”中讨论的原理相同：数据块可以在内存中排序并作为段文件写入磁盘，然后多个排序的段可以合并为一个更大的排序文件。归并排序在磁盘上有良好的顺序访问模式（请记住，优化为顺序I/O是第3章反复出现的主题，这里则再次强调）。

GNU Coreutils (Linux) 中的`sort`实用程序通过自动唤出到磁盘的方式支持处理大于内存的数据集，且排序可以自动并行化以充分利用多核CPU<sup>[9]</sup>。这意味着之前简单的UNIX命令链可以很容易地扩展到大数据集，而不会耗尽内存。从磁盘读取输入文件倒可能会成为瓶颈。

## UNIX设计哲学

我们可以非常容易地使用类似例子中命令链来分析日志文件，这并不是巧合：事实上，这是UNIX的关键设计思想之一，而且时至今日依然令人惊奇。这里，我们更深入地研究一下，目的是从中挖掘、借鉴一些不错的想法<sup>[10]</sup>。







UNIX管道的发明者Doug McIlroy在1964年首先描述了这种用例<sup>[11]</sup>：“当需要以另一种方式处理数据时，我们应该有一些连接程序的方法，就像花园软管互相拧在一起。这就是计算机的I/O。”管道类比一直沿用了下来，通过管道将程序连接起来的想法成为如今的UNIX哲学，在开发人员和UNIX用户中逐渐变得流行的一系列设计原则。1978年，对于这种哲学有了更为完整的描述<sup>[12, 13]</sup>：

1. 每个程序做好一件事。如果要做新的工作，则建立一个全新的程序，而不是通过增加新“特征”使旧程序变得更加复杂。
2. 期待每个程序的输出成为另一个尚未确定的程序的输入。不要将输出与无关信息混淆在一起。避免使用严格的表格状或二进制输入格式。不要使用交互式输入。
3. 尽早尝试设计和构建软件，甚至是操作系统，最好在几周内完成。需要扔掉那些笨拙的部分时不要犹豫，并立即进行重建。
4. 优先使用工具来减轻编程任务，即使你不得不额外花费时间去构建工具，并且预期在使用完成后会将其中一些工具扔掉。

这种方法（自动化、快速原型设计、增量式迭代、测试友好、将大型项目分解为可管理的模块等）听起来非常像今天的敏捷开发与DevOps运动。令人惊讶的是，四十年来并没发生太大变化。

sort命令是一个很好的“一个程序只做好一件事情”的例子。可以说它比大多数编程语言标准库内置的排序算法实现的更好（后者通常不会自动唤出到磁盘，且不使用多线程，即使这样做会更加利于算法）。然而，单独使用sort工具几乎用处不大。只有将其与其他UNIX工具（如uniq）结合使用时，它才会变得无比强大。

像bash这样的UNIX shell可以让我们轻松地将这些小程序组合成强大的数据处理作业。尽管这些程序中是由不同人所编写的，但它们可以灵活地结合在一起。那么，UNIX是如何实现这种可组合性的呢？

## 统一接口

如果希望某个程序的输出成为另一个程序的输入，也就意味着这些程序必须使用相同的数据格式，换句话说，需要兼容的接口。如果你希望能够将任何程序的输出连接到任何程序的输入，那意味着所有程序都必须使用相同的输入/输出接口。

在UNIX中，这个接口就是文件（更准确地说，是文件描述符），文件只是一个有序的字节序列。这是一个非常简单的接口，因此可以使用相同的接口表示许多不同的东西：文件系统上的实际文件，到另一个进程（UNIX socket, stdin, stdout）的通信





通道，设备驱动程序（比如/dev/audio或/dev/lp0），表示TCP连接的套接字等。说起来容易，但实际上对于这些差异明显的不同的事物通过共享统一的接口，使得它们能够轻松连接在一起，确实不可思议<sup>注2</sup>。

按照惯例，许多（但不是全部）UNIX程序将这个字节序列视为ASCII文本。我们的日志分析示例也基于这个事实：awk, sort, uniq和head都将它们的输入文件视为由\n（换行符，ASCII 0x0A）字符分隔的记录列表。其实ASCII记录分隔符0x1E或许是一个更好的选择，因为它本来就是为了这个目的而设计的<sup>[14]</sup>。但是无论如何，所有这些程序都使用标准相同的记录分隔符以支持交互操作。

对每条记录（即一行输入）的解析则没有明确定义。UNIX工具通常使用空格或制表符将行分割成字段，但也使用CSV（逗号分隔）、管道符分隔或其他编码字符。即使像xargs这样相当简单的工具也有六个命令行选项，用于指定如何解析输入。

以ASCII文本作为统一接口的数据格式通常工作得很好，但这种方式并不见得很漂亮：我们的日志分析示例使用{print \$7}来提取网址，但可读性并不好。理想情况下，它看起来应该类似于{print \$request\_url}。我们稍后会再讨论这个想法。

尽管经过几十年的发展依然不够完美，但UNIX统一接口的表现仍然称得上卓越。不过，并没有太多软件能像UNIX工具一样实现交互操作与组合：不能通过自定义分析工具轻松地将电子邮件内容和在线购物历史记录传送到电子表格中，并将结果发布到社交网络或维基百科上。可以说，今天能像UNIX工具一样流畅地连接多个程序绝对是一个“异类”，而并不是常态。

即使是具有相同数据模型的数据库，从一个数据库中导出数据然后导入到另一个数据库也并非易事。集成性的缺失导致了数据的巴尔干化（或称碎片化）。

## 逻辑与布线分离

UNIX工具的另一个特点是使用标准输入（stdin）和标准输出（stdout）。如果运

---

注2：统一接口的另一个例子是URL和HTTP，它们是Web的基础。URL标识网站上的特定事物（资源），并且可以链接到任何其他网站的任何URL。用户通过网络浏览器点击链接在网站之间无缝跳转，即使这些服务器可能是由完全不相关的组织负责的。这个原则现在看来似乎显而易见，但它是网络获得成功的关键。之前的系统并不完全遵守统一原则：例如，在电子公告栏（BBS）时代，每个系统都有自己的电话号码和波特率配置。从一个BBS到另一个BBS的引用说明必须包括电话号码和调制解调器设置，而且用户不得不挂断电话，才能拨打其他BBS，然后手动寻找信息。直接链接到另一个BBS内的某些内容是不可能的。





行一个程序而不指定任何参数，那么标准输入来自键盘，标准输出为屏幕。当然，也可以将文件作为输入和/或将输出重定向到文件。管道允许将一个进程的`stdout`附加到另一个进程的`stdin`（具有小的内存缓冲区，而不需要将全部中间数据流写入磁盘）。

程序仍然可以在需要时直接读取和写入文件。但如果程序不依赖特定的文件路径，只使用`stdin`和`stdout`，则UNIX方法的效果最好。这允许shell用户以任何他们想要的方式连接输入和输出；程序并不知道也不关心输入来自哪里以及输出到哪里。也可以说这是一种松耦合，后期绑定<sup>[15]</sup>或控制反转<sup>[16]</sup>。将输入/输出的布线连接与程序逻辑分开，可以更容易地将小工具组合成更大的系统。

用户甚至可以编写自己的程序，并将它们与操作系统提供的工具组合在一起。程序只需要从`stdin`读取输入并输出至`stdout`，从而参与数据处理流水线。在日志分析示例中，我们可以编写一个工具，将用户代理字、串转换为更为合理的浏览器标识符，或者将IP地址转换为国家代码，并将其插入流水线中。`sort`程序其实并不关心它正在与操作系统进行通信，还是与用户编写的程序进行通信。

但是，使用`stdin`和`stdout`也有其局限性。需要多个输入或输出的程序会变得很棘手。用户不能将程序的输出传输给一个网络连接<sup>[17, 18]注3</sup>。如果一个程序直接打开文件进行读写，或者将另一个程序作为子进程启动，或者打开一个网络连接，那么这个I/O就被程序本身连接起来。虽然支持一些可配置选项（例如通过命令行），但是减少了在shell中连接输入和输出的灵活性。

## 透明与测试

UNIX工具如此成功的部分原因在于，它可以非常轻松地观察事情的进展：

- UNIX命令的输入文件通常被视为是不可变的。这意味着可以随意运行命令，尝试各种命令行选项，而不会损坏输入文件。
- 可以在任何时候结束流水线，将输出管道输送到`less`，然后查看它是否具有预期的形式。这种检查能力对调试非常有用。
- 可以将流水线某个阶段的输出写入文件，并将该文件用作下一阶段的输入。这使得用户可以重新启动后面的阶段，而无需重新运行整个流水线。

---

注3：除非使用一个单独的工具，如`netcat`或`curl`。UNIX试图将所有东西都表示为文件，但是BSD套接字API却偏离了这个惯例<sup>[17]</sup>。研究性操作系统`Plan 9`和`Inferno`在文件使用方面更加一致：它们将TCP连接表示为`/net/tcp`中的文件<sup>[18]</sup>。





因此，与关系数据库的查询优化器相比，尽管UNIX工具相当简单，但是非常有用，特别是对于测试而言。

然而，UNIX工具的最大局限在于它们只能在一台机器上运行，而这正是像Hadoop这样的工具的工作场景。

## MapReduce与分布式文件系统

MapReduce有点像分布在数千台机器上的UNIX工具。与UNIX工具类似，这是一个相当直接、蛮力，却又有效的神奇组件。MapReduce作业可以和UNIX进程相媲美：需要一个或多个输入，并产生一个或多个输出。

和大多数UNIX工具一样，运行MapReduce作业通常不会修改输入，除了生成输出外没有任何副作用。输出文件以序列方式一次性写入（在写文件时，不会修改任何现有的文件）。

UNIX工具使用stdin和stdout作为输入和输出，而MapReduce作业在分布式文件系统上读写文件。在Hadoop的MapReduce实现中，该文件系统被称为HDFS [Hadoop Distributed File System，一个Google文件系统（GFS）的开源实现版本<sup>[19]</sup>]

除HDFS外，还有其他各种分布式文件系统，包括GlusterFS和Quantcast File System（QFS）<sup>[20]</sup>等。诸如Amazon S3，Azure Blob存储和OpenStack Swift<sup>[21]</sup>对象存储服务也有很多相似之处<sup>注4</sup>。在这一章我们主要以HDFS为运行实例，不过这些原则也适用于其他分布式文件系统。

与网络连接存储（NAS）和存储区域网络（SAN）架构的共享磁盘方法相比，HDFS基于无共享原则（参见第二部分的介绍）。共享磁盘存储由集中式存储设备实现，通常使用定制硬件和特殊网络基础设施（如光纤通道）。而无共享方法则不需要特殊硬件，只需要通过传统数据中心网络连接的计算机。

HDFS包含一个在每台机器上运行的守护进程，并会开放一个网络服务以允许其他节点访问存储在该机器上的文件（假设数据中心的每台节点都附带一些本地磁盘）。名为NameNode的中央服务器会跟踪哪个文件块存储在哪家机器上。因此，从概念上

---

注4： 一个不同之处在于，对于HDFS来说，可以将计算任务安排在存储特定文件副本的计算机上运行，而对对象存储通常将存储和计算分开。如果网络带宽是一个瓶颈，从本地磁盘读取文件就会具有性能优势。但是请注意，如果使用纠删编码，将会丢失局部性优势，因为来自多台机器的数据必须进行合并以重建原始文件<sup>[20]</sup>。

讲，HDFS创建了一个庞大的文件系统，来充分利用每个守护进程机器上的磁盘资源。

考虑到机器和磁盘的容错，文件块被复制到多台机器上。复制意味着位于多个机器上的相同数据的多个副本（参阅第5章）；或者像Reed-Solomon代码这样的纠删码方案，相比于副本技术，纠删码可以以更低的存储开销来恢复数据<sup>[20,22]</sup>。这些技术与RAID相似，它可以在同台机器的多个磁盘级别提供数据冗余；不同的是，在分布式文件系统中，文件访问和复制是在传统的数据中心网络上完成的，而不依赖特殊的硬件。

HDFS具有很好的扩展性：在撰写本书时，最大的HDFS集群运行在上万台机器上，总存储容量达到了几百PB [23]。如此大规模的部署主要得益于商用硬件和开源软件的使用，使得HDFS上的数据存储与访问成本远低于同等容量的专用存储设备<sup>[24]</sup>。

## MapReduce作业执行

MapReduce是一个编程框架，可以使用它编写代码来处理HDFS等分布式文件系统的大型数据集。最简单的理解方法是参考本章前面的“简单日志分析”中的Web日志分析示例。MapReduce中的数据处理模式与此非常相似：

1. 读取一组输入文件，并将其分解成记录。在Web日志示例中，每个记录都是日志中的一行（即\n是记录分隔符）。
2. 调用mapper函数从每个输入记录中提取一个键值对。在前面的例子中，mapper函数是`awk '{print $7}'`：它提取URL（\$7）作为关键字，并将相应的值留为空。
3. 按关键字将所有的键值对排序。在日志示例中，这由第一个`sort`命令完成。
4. 调用reducer函数遍历排序后的键值对。如果同一个键出现多次，排序会使它们在列表中相邻，所以很容易组合这些值，而不必在内存中保留过多状态。在前面的例子中，reducer是由`uniq -c`命令实现的，该命令对具有相同关键字的相邻记录进行计数。

这四个步骤可以由一个MapReduce作业执行。步骤2（map）和4（reduce）是用户编写自定义数据处理的代码。步骤1（将文件分解成记录）由输入格式解析器处理。步骤3中的排序步骤`sort`隐含在MapReduce中，无需用户编写，mapper的输出始终会在排序之后再传递给reducer。

要创建MapReduce作业，需要实现两个回调函数，即mapper和reducer，其行为如下（另请参阅第2章的“MapReduce查询”）：

### Mapper

每个输入记录都会调用一次mapper程序，其任务是从输入记录中提取关键字和值。对于每个输入，它可以生成任意数量的键值对（包括空记录）。它不会保留从一个输入记录到下一个记录的任何状态，因此每个记录都是独立处理的。

### Reducer

MapReduce框架使用由mapper生成的键值对，收集属于同一个关键字的所有值，并使用迭代器调用reducer以使用该值的集合。Reducer可以生成输出记录（例如相同URL出现的次数）。

在Web日志示例中，第5步由第二个sort命令按请求数对URL进行排序。在MapReduce中，如果需要第二个排序阶段，则可以编写另一个MapReduce作业并将第一个作业的输出用作第二个作业的输入。这样来看的话，mapper的作用是将数据放入一个适合排序的表单中，而reducer的作用则是处理排序好的数据。

## MapReduce的分布式执行

MapReduce与UNIX命令管道的主要区别在于它可以跨多台机器并行执行计算，其不必编写代码来指示如何并行化。mapper和reducer一次只能处理一条记录，它们不需要知道输入来自哪里或者输出到什么地方，所以框架可以处理复杂的跨机器移动数据的情形。

在分布式计算中可以使用标准的UNIX工具作为mapper和reducer<sup>[25]</sup>，但更为常见的是用传统编程语言实现的函数。在Hadoop MapReduce中，mapper和reducer都是实现特定接口的Java类。在MongoDB和CouchDB中，mapper和reducer是JavaScript函数（参阅第2章的“MapReduce查询”）。

图10-1展示了Hadoop MapReduce作业中的数据流。其并行化基于分区（参阅第6章）实现：作业的输入通常是HDFS中的一个目录，且输入目录中的每个文件或文件块都被视为一个单独的分区，可以由一个单独的map任务来处理（在图10-1标记为map任务1，map任务2和map任务3）。

一个输入文件的大小通常是几百兆字节。只要有足够的空闲内存和CPU资源，MapReduce调度器（图中未显示）会尝试在输入文件副本的某台机器上运行mapper任务<sup>[26]</sup>。这个原理被称为将计算靠近数据<sup>[27]</sup>：它避免将输入文件通过网络进行复制，减少了网络负载，提高了访问局部性。



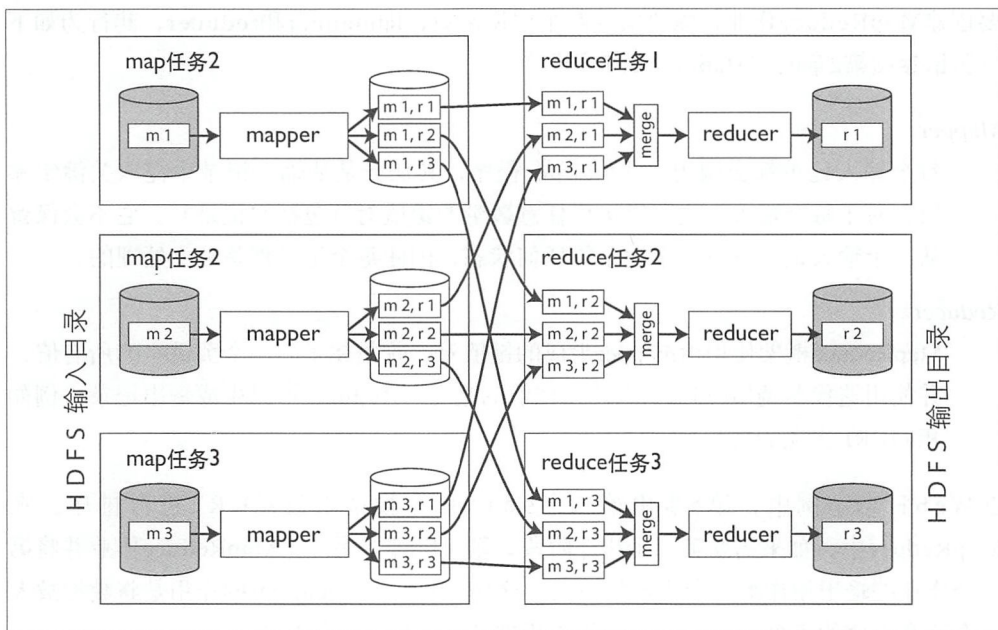


图10-1：带有三个mapper和三个reducer的MapReduce作业

大多数情况下，将要在map任务中所运行的应用程序代码在分配运行任务的节点上并不存在，所以MapReduce框架首先要复制代码（例如Java程序中的JAR文件）到该节点。然后启动map任务并开始读取输入文件，每次将一条记录传递给回调函数mapper。mapper的输出由键值对组成。

Reduce任务中的计算也被分割成块。Map任务的数量由输入文件块的数量决定，而reduce任务的数量则是由作业的作者来配置的（可以不同于map任务的数量）。为了确保具有相同关键字的所有键值对都在相同的reducer任务中处理，框架使用关键字的哈希值来确定哪个reduce任务接收特定的键值对（请参阅第6章“基于关键字哈希值分区”）。

键值对必须进行排序。如果数据集太大，可能无法在单台机器上使用常规排序算法。事实上，排序是分阶段进行的。首先，每个map任务都基于关键字哈希值，按照reducer对输出进行分块。每一个分区都被写入mapper程序所在本地磁盘上的已排序文件，使用的技术类似我们在第3章讨论的“SSTables和LSM-Trees”。

当mapper完成读取输入文件并写入经过排序的输出文件，MapReduce调度器就会通知reducer开始从mapper中获取输出文件。reducer与每个mapper相连接，并按照其分区从mapper中下载排序后的键值对文件。按照reducer分区，排序和将数据分区从mapper

复制到reducer，这样一个过程被称为*shuffle*<sup>[26]</sup>（这是一个容易令人困惑的术语，它并不完全与洗牌一样，在MapReduce中其实没有随机性）。

reduce任务从mapper中获取文件并将它们合并在一起，同时保持数据的排序。因此，如果不同的mapper使用相同的关键字生成记录，则这些记录会在合并后的reducer输入中位于相邻位置。

reducer通过关键字和迭代器进行调用，而迭代器逐步扫描所有具有相同关键字的记录（某些情况下可能无法完全在内存中完成）。reducer可以使用任意逻辑来处理这些记录，并且生成任意数量的输出记录。这些输出记录被写入分布式文件系统中的文件（通常是在运行reducer机器的本地磁盘中的一个拷贝，在其他机器上存在副本）。

## MapReduce工作流

单个MapReduce作业可以解决的问题范围有限。回顾一下日志分析的例子，一个MapReduce作业可以确定每个URL页面的浏览次数，但不是最受欢迎的那些URL，因为这需要第二轮排序。

因此，将MapReduce作业链接到工作流中是非常普遍的，这样，作业的输出将成为下一个作业的输入。Hadoop MapReduce框架对工作流并没有任何特殊的支持，所以链接方式是通过目录名隐式完成的：第一个作业必须配置为将其输出写入HDFS中的指定目录，而第二个作业必须配置为读取相同的目录名作为输入。从MapReduce框架的角度来看，它们仍然是两个独立的作业。

因此，链接方式的MapReduce作业并不像UNIX命令流水线（它直接将一个进程的输出作为输入传递至另一个进程，只需要很小的内存缓冲区），而更像是一系列命令，其中每个命令的输出被写入临时文件，下一个命令从临时文件中读取。这种设计有利有弊，我们将在本章后面的“中间状态实体化”中讨论。

只有当作业成功完成时，批处理作业的输出才会被视为有效（MapReduce会丢弃失败作业的部分输出）。因此，工作流中的一个作业只有在先前的作业（即生成其输入目录的作业）成功完成时才能开始。为了处理这些作业执行之间的依赖关系，已经开发了各种Hadoop的工作流调度器，包括Oozie, Azkaban, Luigi, Airflow和Pinball<sup>[28]</sup>。

这些调度程序还具有管理功能，在维护大量批处理作业时非常有用。在构建推荐系统<sup>[29]</sup>时，由50到100个MapReduce作业组成的工作流是非常常见的。而在大型组织中，许多不同的团队可能运行不同的作业来读取彼此的输出。良好的工具支持对于管理如此复杂的数据流非常重要。

Hadoop的各种高级工具（如Pig<sup>[30]</sup>，Hive<sup>[31]</sup>，Cascading<sup>[32]</sup>，Crunch<sup>[33]</sup>和FlumeJava<sup>[34]</sup>）则支持设置多个MapReduce阶段的工作流，这些不同的阶段会被恰当地自动链接起来。

## Reduce端的join与分组

我们在第2章中讨论了数据模型和查询语言的联结操作，但是我们还没有深入探讨join是如何实现的。现在我们将再次开始这个话题。

在许多数据集中，通常一条记录会与另一条记录存在关联，例如：关系模型中的外键，文档模型中的文档引用或图模型中的边。只要有代码需要访问该关联两边的记录（包含引用的记录和被引用的记录），那么就需要join操作。正如第2章所讨论的那样，反规范化可以减少对join的需求，但通常无法完全避免join<sup>注5</sup>。

在数据库中，如果执行的查询只涉及少量记录，那么数据库通常会使用索引来加速查找（请参阅第3章）。如果查询涉及到join操作，则可能需要对多个索引进行查找。然而，MapReduce没有索引的概念，至少不是通常意义上的索引。

当给定一组文件作为MapReduce输入时，它读取所有文件的全部内容；数据库将其称作全表扫描。如果只想读取少量记录，则与索引查找相比，全表扫描的成本非常昂贵。但是，分析查询（参阅第3章“事务处理与分析处理”）通常需要计算大量记录的聚合。在这种情况下，扫描整个数据集是比较合理的，特别是如果可以在多台机器上并行处理。

在批处理的背景下讨论join时，我们主要是解决数据集内存在关联的所有事件。例如，假设一个作业是同时为所有用户处理数据，而不仅仅是为一个特定的用户查找数据（这可以通过索引更高效地完成）。

### 示例：分析用户活动事件

图10-2给出了批处理作业中典型的join示例。图中左侧是事件日志，描述登录用户在网站上的活动（称为活动事件或单击流数据），右侧是用户数据库。可以将此示例视为一种星型模式的一部分（请参阅第3章的“星型与雪花型分析模式”）：事件日志是事实数据表，用户数据库是维度之一。

---

注5：我们在本书中讨论的联结通常是等值联结，一种最常见的联结类型，是指记录与在特定字段（如ID）中具有相同值的其他记录相关联。某些数据库支持更通用的联结类型，例如使用小于运算符而不是等号运算符，但这不在本书的讨论范围之内。



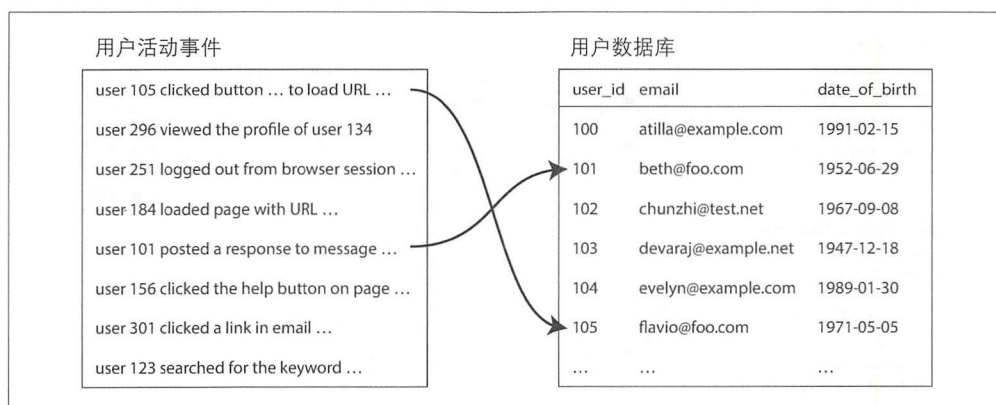


图10-2：用户活动事件日志与用户描述数据库之间的join

分析任务可能需要将用户活动与用户描述信息相关联：例如，如果描述中包含用户年龄或出生日期，则系统可以确定哪些年龄组最受欢迎。但是，活动事件中仅包含用户标识，而不包含完整的用户描述信息。而在每个活动事件中嵌入这些描述信息又会太浪费。因此，活动事件需要与用户描述数据库进行join。

join的最简单实现是逐个遍历活动事件，并在（远程服务器上的）用户数据库中查询每个遇到的用户ID。该方案首先是可行的，但性能会非常差：吞吐量将受到数据库服务器的往返时间的限制，本地缓存的有效性将很大程度上取决于数据的分布，并且同时运行的大量并行查询很容易使数据库不堪重负<sup>[35]</sup>。

为了在批处理过程中实现良好的吞吐量，计算必须（尽可能）在一台机器上进行。如果通过网络对每条记录进行随机访问则请求太慢。而且，考虑到远程数据库中的数据可能会发生变化，查询远程数据库意味着会增加批处理作业的不确定性。

因此，更好的方法是获取用户数据库的副本（例如，使用ETL进程从数据库备份中提取数据，请参阅第3章的“数据仓库”），并将其放入与用户活动事件日志相同的分布式文件系统。然后，可以将用户数据库放在HDFS中的一组文件中，并将用户活动记录放在另一组文件中，使用MapReduce将所有相关记录集中到一起，从而有效地处理它们。

## 排序-合并join

回想一下，mapper的目的是从每个输入记录中提取关键字和值。在图10-2的情况下，这个关键字就是用户ID：一组mapper会扫描活动事件（提取用户ID作为关键字，而活动事件作为值），而另一组mapper将会遍历用户数据库（提取用户ID作为关键字，用户出生日期作为值）。过程如图10-3所示。

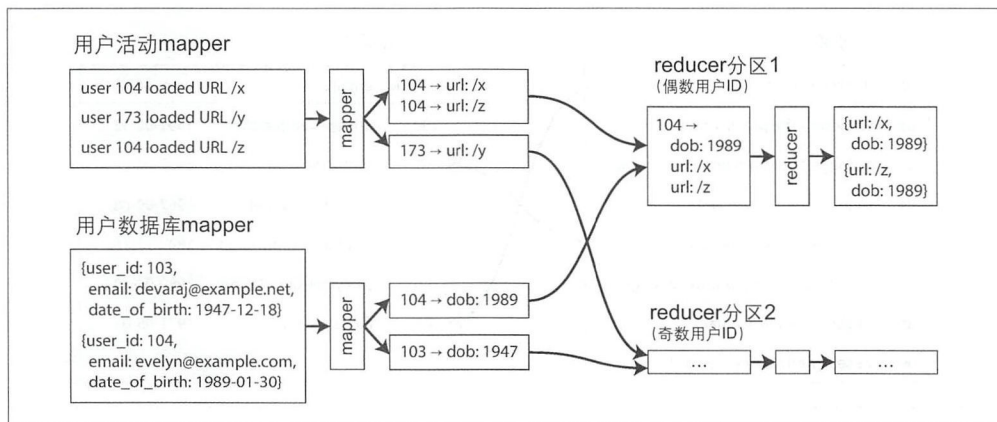


图10-3：用户ID的reduce端排序-合并join。如果输入数据集被分割成多个文件，则每个文件可以被多个mapper并行处理

当MapReduce框架通过关键字对mapper输出进行分区，然后对键值对进行排序时，结果是所有活动事件和用户ID相同的用户记录在reducer的输入中彼此相邻。MapReduce作业甚至可以对记录进行排序，以便reducer会首先看到用户数据库中的记录，然后按时间戳顺序查看活动事件，这种技术称为次级排序<sup>[26]</sup>。

然后reducer可以很容易地执行真正的join逻辑：为每个用户ID调用一次reducer函数。由于次级排序，第一个值应该是来自用户数据库的出生日期记录。Reducer将出生日期存储在局部变量中，然后使用相同的用户ID遍历活动事件，输出相应的已观看网址和观看者年龄。随后的MapReduce作业可以计算每个URL的查看者年龄分布，并按年龄组进行聚类。

由于reducer每次处理一个特定用户ID的所有记录，因此只需要将用户记录在内存中保存一次，而不需要通过网络发出任何请求。这个算法被称为排序-合并join，因为mapper的输出是按关键字排序的，然后reducer将来自join两侧的已排序记录列表合并在一起。

## 将相关数据放在一起

在排序-合并join中，mapper和排序过程确保将执行特定用户IDjoin操作的所有必要数据都放在一起，这样就只需要一次reducer调用。因为所有需要的数据已经预先排列好，所以reducer是一段相当简单的单线程代码，以高吞吐量和低内存开销来处理记录。



一种理解这种架构的方法是mapper发送“消息”给reducer。当mapper发出一个键值对时，关键字就像传递值的目标地址一样。即使关键字只是一个任意字符串（不是那种像IP地址和端口一样的实际网络地址），但是它的行为就像一个地址：所有具有相同关键字的键值对都将被传送到相同的目的地（即对reducer的调用）。

使用MapReduce编程模型将计算中的物理网络通信部分（从正确的机器获取数据）从应用逻辑（处理数据）中分离出来。这种分离与数据库的典型使用方式形成了鲜明对比：从数据库中获取数据的请求经常发生在应用程序代码的深处<sup>[36]</sup>。由于MapReduce能够处理所有的网络通信，因此它也避免了在应用程序代码中处理局部故障，例如某个节点的崩溃：MapReduce会在不影响应用程序逻辑的情况下透明地重试失败的任务。

## 分组

除了join之外，“将相关数据放在一起”模式的另一个常见用法是通过某个关键字（如SQL中的GROUP BY子句）对记录进行分组。所有具有相同关键字的记录形成一个组，然后在每个组内执行某种聚合操作，例如：

- 计算每个组中记录的数量 [例如，在统计页面视图的示例中，SQL将其表示为COUNT (\*) 聚合] 。
- 对SQL中的特定字段进行求和 [SUM(fieldname)] 。
- 根据排名函数选择前k个记录。

使用MapReduce实现这种分组操作的最简单方法是设置mapper，使其生成的键值对使用所需的分组关键字。然后，分区和排序过程将相同reducer中所有具有相同关键字的记录集合在一起。因此，在MapReduce上实现的分组和join看起来非常相似。

分组的另一个常见用途是收集特定用户会话的所有活动事件，以便发现用户的活动序列，称为会话流程<sup>[37]</sup>。例如，可以使用这种分析来确定选择网站新版本的用户是否比选择旧版本（A / B测试）的用户更有可能产生购买行为，或计算某个营销活动是否有效。

如果有多个Web服务器处理用户请求，则特定用户的活动事件很可能分散在各个不同服务器的日志文件中。可以通过使用会话cookie、用户ID或类似的标识符作为分组关键字来实现访问流程，将特定用户的所有活动事件放在一起，同时将不同用户的事件分配到不同的分区。





## 处理数据倾斜

如果与单个关键字相关的数据量非常大，那么会破坏掉“将所有具有相同关键字的记录放在一起”的模式。例如，在社交网络中，大多数用户会有上百人关注者，但少数名人则可能有数百万的追随者。这种不成比例的活跃数据库记录被称为关键对象<sup>[38]</sup>或热键。

在单个reducer中收集与名人相关的所有活动（例如对他们发布内容的回复）可能会导致严重的数据倾斜（也称为热点）。也就是说，某个reducer必须处理比其他reducer更多的记录（请参阅第6章的“负载倾斜与热点”）。由于MapReduce作业只有在其所有mapper和reducer都完成时才能完成，因此所有后续作业必须等待最慢的reducer完成之后才能开始。

如果join输入中存在热键，则可以使用算法进行补偿。例如，Pig中的倾斜join方法首先运行一个抽样作业来确定哪些属于热键<sup>[39]</sup>。在真正开始执行join时，mapper将任何与热键有关的记录发送到随机选择的若干个reducer中的一个（传统MapReduce基于关键字哈希来确定性地选择reducer）。对于join的其他输入，与热键相关的记录需要被复制到所有处理该关键字的reducer中<sup>[40]</sup>。

这种技术将处理热键的工作分散到多个reducer上，可以更好地实现并行处理，代价是不得不将join的其他输入复制到多个reducer。Crunch中的共享join方法与此类似，但需要明确指定热键，而不是使用抽样作业。这种技术也非常类似我们在第6章“负载倾斜与热点”所讨论的技术，使用随机化来缓解分区数据库中的热点。

Hive的倾斜join优化采取了另一种方法。它需要在表格元数据中明确指定热键，并将与这些键相关的记录与其余文件分开存放。在该表上执行join时，它将对热键使用map端join（请参阅下一节）。

使用热键对记录进行分组并汇总时，可以分两个阶段进行分组。第一个MapReduce阶段将记录随机发送到reducer，以便每个reducer对热键的记录子集执行分组，并为每个键输出更紧凑的聚合值。然后第二个MapReduce作业将来自所有第一阶段reducer的值合并为每个键的单一值。

## map端join操作

上一节描述的join算法在reducer中执行实际的join逻辑，因此被称为reduce端join。mapper负责准备输入数据：从每个输入记录中提取关键字和值，将键值对分配给reducer分区，并按关键字排序。



Reduce端join方法的优点是不需要对输入数据做任何假设：无论其属性与结构如何，mapper都可以将数据处理好以准备join。然而，不利的一面是，所有这些排序，复制到reducer以及合并reducer输入可能会是非常昂贵的操作，这取决于可用的内存缓冲区，当数据通过MapReduce阶段时，数据可能需要写入磁盘若干次<sup>[37]</sup>。

另一方面，如果可以对输入数据进行某些假设，则可以通过使用所谓的map端join来加快速度。这种方法使用了一个缩减版本的MapReduce作业，其中没有reducer，也没有排序；相反，每个mapper只需从分布式文件系统中读取输入文件块，然后将输出文件写入文件系统即可。

## 广播哈希join

实现map端join的最简单方法特别适合大数据集与小数据集join，尤其是小数据集能够全部加载到每个mapper的内存中。

例如，假设对于图10-2情况，用户数据库可以完全放入内存。在这种情况下，当mapper程序执行时，它可以首先将用户数据库从分布式文件系统读取到内存的哈希表中。然后，mapper程序扫描用户活动事件，并简单地查找哈希表中每个事件的用户ID<sup>注6</sup>。

Map任务依然可以有多个：大数据集的每个文件块对应一个mapper（在图10-2的例子中，活动事件是大输入数据集）。每个mapper还负责将小数据集全部加载到内存中。

这种简单而有效的算法被称为广播哈希join：“广播”一词主要是指大数据集每个分区的mapper还读取整个小数据集（即小数据集实际被“广播”给大数据集），“哈希”意味着使用哈希表。有多种系统都支持该方法，包括Pig（称为replicated join），Hive（称为MapJoin），Cascading和Crunch等。它也用于数据仓库查询引擎，例如Impala<sup>[41]</sup>。

另一种方法并不需要将小数据集加载至内存哈希表中，而是将其保存在本地磁盘上的只读索引中<sup>[42]</sup>。由于频繁访问，索引大部分内容其实是驻留在操作系统的页面缓存中，因此这种方法可以提供与内存哈希表几乎一样快的随机访问性能，而实际上并不要求整个数据集读入内存。

---

注6： 这个例子假定哈希表中的每个键只有一个条目，这对用户数据库（ID唯一标识一个用户）来说往往是正确的。通常，哈希表可能包含具有相同键的多个条目，并且联结运算符会输出关键字的所有匹配记录。





## 分区哈希join

如果以相同方式对map端join的输入进行分区，则哈希join方法可以独立作用于每个分区。在图10-2的情况下，可以根据用户ID的最后一位十进制数字（因此每一边都有10个分区）来分配活动事件和用户数据库中的记录。例如，Mapper 3首先将所有以3结尾的ID的用户加载到哈希表中，然后扫描ID以3结尾的每个用户的所有活动事件。

如果分区操作正确完成，就可以确定所有要join的记录都位于相同编号的分区中，因此每个mapper只需从每个输入数据集中读取一个分区就足够了。这样的优点是每个mapper都可以将较少的数据加载到其哈希表中。

这种方法只适用于两个join的输入具有相同数量的分区，根据相同的关键字和相同的哈希函数将记录分配至分区。如果输入是由之前已经执行过这个分组的MapReduce作业生成的，那么这是一个合理的假设。

分区哈希join在Hive<sup>[37]</sup>中称为*bucketed map join*。

## map端合并join

如果输入数据集不仅以相同的方式进行分区，而且还基于相同的关键字进行了排序，则可以应用map端join的另一种变体。这时，输入是否足够小以载入内存并不重要，因为mapper可以执行通常由reducer执行的合并操作：按关键字升序增量读取两个输入文件，并且匹配具有相同关键字的记录。

如果map端合并join是可能的，则意味着先前的MapReduce作业会首先将输入数据集引入到这个经过分区和排序的表单中。原则上，join可以在之前作业的reduce阶段进行。但是，在独立的map作业中执行合并join更为合适，例如，除了特定的join操作之外，分区和排序后的数据集还可用于其他目的。

## 具有map端join的MapReduce工作流

当下游作业使用MapReduce join的输出时，map端或reduce端join的不同选择会影响到输出结构。reduce端join的输出按join关键字进行分区和排序，而map端join的输出按照与大数据集相同的方式进行分区和排序（因为对大数据集的每个文件块都会启动一个map任务，无论是使用分区join还是广播join）。

正如所讨论的，map端join也存在对输入数据集的大小、排序和分区方面的假设。在优化join策略时，了解分布式文件系统中数据集的物理布局非常重要：仅仅知道编码格式和数据存储目录的名称是不够的；还必须知道数据分区数量，以及分区和排序的关键字。







在Hadoop生态系统中，关于数据集分区的元数据经常在HCatalog和Hive metastore中维护<sup>[37]</sup>。

## 批处理工作流的输出

我们已经讨论了很多关于实现MapReduce工作流的算法，但仍然忽略了一个重要问题：一旦作业完成，处理的最终结果是什么？为什么要把这些作业放在首位？

数据库查询中，根据分析目的我们区分了事务处理（OLTP）和分析型处理（参阅第3章“事务处理与分析处理”）。我们看到OLTP查询通常使用索引按关键字查找少量记录，然后将查询结果呈现给用户（例如在网页上）。另一方面，分析查询通常会扫描大量记录，执行分组与聚合，输出完整的报告，例如：显示度量随时间变化的图表，或某种排名中的前10项，或将某一数量分解成若干子类别。这种报告的消费者通常是需要做出商业决策的分析师或经理。

批处理即不是事务处理，也不是分析，那么，把它放在哪里更合适？批处理与分析更为接近，因为批处理过程通常会扫描大部分的输入数据集。但是，MapReduce作业的工作流与分析中SQL查询不同（请参阅本章后面的“对比Hadoop与分布式数据库”）。批处理过程的输出通常不是报告，而是其他类型的数据结构。

## 生成搜索索引

Google最初使用MapReduce的目的是为其搜索引擎建立索引，这个索引被实现为5到10个MapReduce作业的工作流<sup>[1]</sup>。尽管Google后来不再使用MapReduce<sup>[43]</sup>，但是如果从构建搜索索引的角度来看MapReduce，会更加有助于理解MapReduce（即使在今天，Hadoop MapReduce仍然是构建Lucene/Solr索引的好方法<sup>[44]</sup>）。

在第3章的“全文搜索与模糊索引”中简要讨论了像Lucene这样的全文搜索索引是如何工作的：它是一个文件（术语字典），可以在其中有效地查找特定关键字，并找到包含该关键字的所有文档ID列表（发布列表）。这是一个非常简单的搜索索引视图，实际上它还需要各种附加数据，以便按相关性对检索结果进行排序、拼写检查、解析同义词等等，但基本原则类似。

如果需要对一组固定文档进行全文检索，则批处理是构建索引的有效方法：mapper根据需要对文档集进行分区，每个reducer构建其分区索引，并将索引文件写入分布式文件系统。并行处理非常适用于构建这样的文档分区索引（请参阅第5章的“分区与二级索引”）。

由于按关键字查询搜索索引是只读操作，因此这些索引文件一旦创建就是不可变的。





如果索引的文档集合发生更改，则可以选择定期重新运行整个索引工作流，并在完成后用新的索引文件批量替换之前的索引文件。如果只有少量文档发生了变化，这种方法在计算上可能比较昂贵，但是它的优点是索引过程非常清晰合理：文档作为输入，索引作为输出。

增量建立索引是一种替代方法。如第3章所述，如果要添加、删除或更新索引中的文档，Lucene会生成新的段文件，并在后台异步合并和压缩段文件。我们将在第11章中看到更多这样的增量处理。

## 批处理输出键值

搜索索引只是批处理工作流输出的一个示例。批处理的另一个常见用途是构建机器学习系统，如分类器（例如垃圾邮件过滤器，异常检测，图像识别）和推荐系统（例如你可能认识的人，你可能感兴趣的产品或相关搜索<sup>[29]</sup>）。

这些批量作业的输出通常是某种数据库：例如，在用户数据库中通过用户ID进行查询以获取建议的好友，或者在产品数据库中通过产品ID查询以获取相关产品列表<sup>[45]</sup>。

查询数据库需要在处理用户请求的Web应用中进行，而这些请求通常与Hadoop基础架构是分离的。那么批处理过程的输出如何返回至数据库中以供Web应用查询？

最明显的选择可能是直接在mapper或reducer中使用你最喜欢的数据库客户端软件包，而批处理作业则直接写入至数据库服务器，一次写入一条记录。这样的方法可行（假设防火墙规则允许从Hadoop环境直接访问生产数据库），但由于以下原因，这并不是一个好方案：

- 正如前面讨论join时提到的，为每个记录发送一个网络请求比批处理任务的正常吞吐量要慢几个数量级。即使客户端软件包支持批处理，性能也可能很差。
- MapReduce作业经常并行处理许多任务。如果所有的mapper或reducer都同时写入同一个输出数据库，并以批处理期望的速率写入，那么数据库很容易过载，其查询性能会受到影响。这可能会进一步导致系统其他部分的操作问题<sup>[35]</sup>。
- 通常情况下，MapReduce为作业输出提供了一个干净的“全有或全无”的保证：如果作业成功，则结果就是只运行一次任务的输出，即使中间发生了某些任务失败但最终重试成功。如果整个作业失败，则不会产生输出。然而，从作业内部写入外部系统会产生外部可见的副作用，而这种副作用无法彻底屏蔽。因此，将不得不担心部分完成的作业产生对其他系统可见的结果，以及Hadoop任务尝试和推测性执行的复杂性。





更好的解决方案是，批处理作业创建一个全新的数据库，并将其作为文件写入分布式文件系统上的作业输出目录，就像上一节的搜索索引一样。这些数据文件一旦写入就是不可变的，可以批量加载到处理只读查询的服务器中。各种键值存储都支持构建MapReduce作业中的数据库文件，包括Voldemort<sup>[46]</sup>，Terrapin<sup>[47]</sup>，ElephantDB<sup>[48]</sup>和HBase批量加载<sup>[49]</sup>。

构建数据库文件是很好的MapReduce使用方法：使用mapper提取关键字，然后使用该关键字进行排序，这已经成为构建索引的必要步骤。由于大多数键值存储是只读的（文件只能由批处理作业一次写入，而且是不可变的），所以数据结构非常简单，例如，它们不需要预写日志（请参阅第3章的“可靠的B-tree”）。

将数据加载到Voldemort时，服务器将继续向旧数据文件发起请求，同时将新数据文件从分布式文件系统复制到服务器的本地磁盘。一旦复制完成，服务器会自动切换到新文件进行查询。如果在这个过程中出现任何问题，它可以很容易地再次切换回旧文件，因为这些旧文件依然存在，而且不可变<sup>[46]</sup>。

## 批处理输出的哲学

本章前面讨论过的UNIX设计哲学倡导明确的数据流：一个程序读取输入并写回输出。在这个过程中，输入保持不变，任何以前的输出都被新输出完全替换，并且没有其他副作用。这意味着可以随心所欲地重新运行一个命令，进行调整或调试，而不会扰乱系统状态。

MapReduce作业的输出处理遵循相同的原理。将输入视为不可变，避免副作用（例如对外部数据库的写入），批处理作业不仅实现了良好的性能，而且更容易维护：

- 如果在代码中引入了漏洞，输出错误或者损坏，那么可以简单地回滚到先前版本，然后重新运行该作业，将再次生成正确的输出；或者更简单的办法是将旧的输出保存在不同的目录中，然后切换回原来的目录。相比之下，具有读写事务的数据库就不具有这样的属性：如果部署了将错误数据写入数据库的错误代码，回滚代码并不能修复数据库中的数据。能够从错误代码中恢复的思想被称为人为容错性<sup>[50]</sup>。
- 与发生错误即意味着不可挽回的损害相比，易于回滚的特性更有利于快速开发新功能。这种使不可逆性最小化的原则对于敏捷开发是有益的<sup>[51]</sup>。
- 如果map或reduce任务失败，MapReduce框架会自动重新安排作业并在同一个输入上再次运行。如果失败是由于代码漏洞造成的，那么它会一直崩溃，最终导致作业在数次尝试之后失败。但是如果故障是由于暂时问题引起的，则可以实现





容错。由于输入总是不可变，这种自动重试是安全的，而失败任务的输出则被MapReduce框架丢弃。

- 相同的文件可用作各种不同作业的输入，其中包括监控作业，它可以搜集相关运行指标，并评估其他作业的输出是否满足预期特性（例如，将其与前一次运行的输出进行比较并测量差异）。
- 与UNIX工具类似，MapReduce作业将逻辑与连线（配置输入和输出目录）分开，从而可以更好地隔离问题，重用代码：一个团队可以专注于实现“做好一件事”的工作，而其他团队可以决定何时何地运行该作业。

在这些方面，那些UNIX验证过的设计原则似乎用在Hadoop上也同样效果很好，但UNIX和Hadoop在某些方面存在不同。例如，因为大多数UNIX工具都假定输入为没有类型的文本文件，所以必须做大量的输入解析工作（本章开头的日志分析示例使用`{print $7}`来提取URL）。在Hadoop中，通过使用更多结构化的文件格式，可以消除一些低价值的语法转换：Avro（请参阅第4章的“Avro”）和Parquet（请参阅第3章的“列式存储”）通常用于提供高效的基于模式的编码，并支持模式的不断演变（参见第4章）。

## 对比Hadoop与分布式数据库

正如我们所看到的，Hadoop有点像UNIX的分布式版本，其中HDFS是文件系统，而MapReduce则是UNIX进程的特殊实现（它总是在map阶段和reduce阶段之间运行sort工具）。我们讨论了如何在这些原语之上实现各种join和分组操作。

MapReduce论文<sup>[1]</sup>发表时，在某种意义上来说它并不新鲜。我们在前几节讨论的所有处理和并行join算法都已经在十多年前所谓的大规模并行处理（MPP）数据库中实现了<sup>[3,4]</sup>。Gamma数据库机器，Teradata和Tandem NonStop SQL是这方面的先驱<sup>[52]</sup>。

这其中最大的区别在于MPP数据库专注于在一个机器集群上并行执行SQL查询分析，而MapReduce和分布式文件系统<sup>[19]</sup>的结合则更像一个可以运行任意程序的通用操作系统。

### 存储多样性

数据库要求根据特定的模型（例如，关系或文档）来构造数据，而分布式文件系统中的文件只是字节序列，可以使用任何数据模型和编码来编写。它们可能是数据库记录的集合，也同样可以是文本、图像、视频、传感器读数、稀疏矩阵、特征向量、基因组序列或任何其他类型的数据。



更直接地讲，Hadoop开放了将数据不加区分地转存到HDFS的可能性，之后再去考虑如何进一步处理数据<sup>[53]</sup>。相比之下，在将数据导入数据库的专有存储格式之前，MPP数据库通常需要对数据和查询模式进行仔细的前期建模。

从纯粹主义者的角度来看，这种仔细的建模和导入方式似乎是可取的，因为这意味着数据库用户拥有质量更好的数据。然而在实践中，看起来只是简单地使数据可用；而即使是一个古怪的，难以使用的原始格式往往也比预先决定理想的数据模型更有价值<sup>[54]</sup>。

这个想法与数据仓库类似（请参阅第3章的“数据仓库”）：只将大型组织中各个部分的数据集中在一起是非常有价值的，因为它可以在之前完全分离的数据集执行join操作。MPP数据库所秉持的谨慎设计模式减慢了集中式数据的收集速度；因此，仅仅以原始形式收集数据，之后再考虑模式设计，从而使收集数据的速度加快（这种形式有时也被称为“数据湖”或“企业数据中心”<sup>[55]</sup>）。

不加区分地数据转储也转移了数据解释的负担：不是强迫数据集的生产者将其转化为标准化格式，而是将解释数据变为消费者的问题（读时模式<sup>[56]</sup>，请参阅第2章的“文档模型中的模式灵活性”）。如果生产者和消费者分属不同优先级的团队，这会是一个优势。甚至也可能并不存在一个理想的数据模型，而是针对不同的目的产生对数据不同的看法。以原始形式简单地转储数据可以进行多次这样的转换。这种方法被称为寿司原则<sup>[57]</sup>。

因此，Hadoop经常被用于实现ETL过程（参见第3章“数据仓库”）：来自事务处理系统的数据以某种原始形式转储到分布式文件系统中，然后编写MapReduce作业进行数据清理，将其转换为关系表单，并将其导入MPP数据仓库以进行分析。数据建模仍然会发生，但它位于一个单独步骤中，与数据收集是分离的。由于分布式文件系统支持以任何格式编码的数据，所以这种解耦是可行的。

## 处理模型的多样性

MPP数据库属于一体化、紧密集成的软件系统，包括磁盘存储布局、查询计划、调度和执行。由于这些组件都可以针对数据库的特定需求进行调整和优化，因此整个系统可以在其设计的查询类型上实现非常好的性能。此外，SQL查询语言支持表达式查询，以及优雅的语义，而无需编写代码，因此，业务分析人员使用的图形化工具（如Tableau）即可执行查询。

另一方面，并非所有类型的处理都可以合理地表达为SQL查询。例如，如果正在构建机器学习和推荐系统，或具有相关性排名模型的全文搜索索引，或执行图像分析，则



很可能需要更具一般性的数据处理模型。这些类型的处理通常特定于专门的应用程序（例如机器学习的特征工程，机器翻译的自然语言模型，欺诈预测的风险评估功能），所以不可避免地需要编写代码，而不仅仅是查询。

MapReduce使工程师能够轻松地在大数据集上运行自己的代码。如果你有HDFS和MapReduce，可以在它上面建立一个SQL查询执行引擎，事实上这就是Hive项目所做事情<sup>[31]</sup>。也可以编写许多其他形式的不适合用SQL查询表示的批处理。

随后，人们发现MapReduce对于某些类型的处理局限性太大，而且执行得太差，因此在Hadoop之上开发了各种其他的处理模型（我们将在本章后面的“超越MapReduce”中展开讨论）。仅有SQL和MapReduce两种处理模型还不够，我们需要更多不同的模型！而且由于Hadoop平台的开放性，实施一整套处理方法也是可行的，但对于一体化的MPP数据库来说则力所不及<sup>[58]</sup>。

至关重要的是，这些不同的处理模型都可以在一个共享的机器集群中运行，所有机器都可以访问分布式文件系统上的相同文件。在Hadoop方法中，不需要将数据导入到几个不同的专用系统中进行不同类型的处理：系统已经足够灵活，可以支持同一集群中不同的工作负载。无需移动数据使得从数据中获取价值变得容易得多，而且更容易使用新的处理模型进行测试。

Hadoop生态系统包括可以随机访问的OLTP数据库，如HBase（请参阅第3章的“SSTables和LSM-Tree”）以及MPP模式的分析数据库，如Impala<sup>[41]</sup>。HBase和Impala都不使用MapReduce，但都使用HDFS进行存储。尽管它们访问和处理数据的方法差异很大，但是可以共存并被集成到同一个系统中。

## 针对频繁故障的设计

在比较MapReduce和MPP数据库时，设计方法的另外两个不同点是：如何处理故障和如何使用内存和磁盘。与在线系统相比，批处理对故障的敏感度较低，因为如果遇到失败的任务，它们不会立即影响用户，而且总是可以重新运行。

如果一个节点在执行查询时崩溃，大多数MPP数据库会中止整个查询，并让用户重新提交查询或自动重新运行查询<sup>[3]</sup>。由于查询通常最多运行几秒钟或几分钟，所以这种处理错误的方式是可以接受的，因为重试的代价不是太大。MPP数据库还倾向于在内存中保留尽可能多的数据（例如使用哈希join）以避免磁盘读取的成本。

另一方面，MapReduce可以容忍map或reduce任务的失败，通过以单个任务的粒度重试工作来避免影响整体作业。它也会期望将数据写入磁盘，一方面是为了容错，另一



方面是假设数据集太大而不能载入内存。

MapReduce方法更适用于较大的作业：处理大量的数据并运行很长时间的作业，以至于在此过程中可能至少会遇到一次任务故障。在这种情况下，由于单个任务失败而重新运行整个工作将是巨大的浪费。虽然以单个任务粒度进行恢复带来的开销会使无故障处理变得更慢，但如果任务失败率足够高，这仍然是一种合理的权衡方式。

但是这些假设的现实性究竟如何呢？在大多数集群中，机器故障确实在发生，但并不是很频繁，可能很少发生，大多数作业都不会遇到机器故障。为了容错，真的值得引入这么大的开销吗？

要了解MapReduce节省使用内存和任务级恢复的背后原因，查看最初设计MapReduce的环境会很有帮助。Google拥有混合使用的数据中心，在线生产服务和离线批处理作业在同一台机器上运行。任务通过使用容器来实现资源分配（CPU，内存，磁盘空间等）。每个任务都具有优先级，如果优先级较高的任务需要更多的资源，则可以终止（抢占）同一机器上较低优先级的任务以释放资源。优先级还决定了计算资源的定价：团队必须为他们使用的资源付费，而优先级更高的任务花费也会更多<sup>[59]</sup>。

这种架构允许非生产（低优先级）计算资源被过度使用，因为系统知道可以在必要时回收资源。与那些隔离生产和非生产任务的系统相比，过度使用资源反过来可以更好地利用机器和提高效率。但是，由于MapReduce作业以低优先级运行，当优先级较高的进程需要其资源时，随时都可以发生抢占。批处理任何工作可以有效地利用高优先级进程剩下的任何可用计算资源。

在Google，运行一个小时的MapReduce任务大约有5%被终止的风险，为更高优先级的进程腾出空间。这个比例要比由于硬件问题、机器重启或其他原因导致的故障率高出一个数量级<sup>[59]</sup>。以这样的抢占率，如果一个作业有100个任务，每个任务运行10分钟，那么至少有一个任务在完成之前将被终止的风险大于50%。

这就是为什么MapReduce被设计为容忍意外任务终止的原因：不是因为硬件特别不可靠，而是因为任意终止进程的灵活性能更好地利用集群资源。

而在开源集群调度器中，抢占的使用情况则相对较少。YARN的CapacityScheduler支持抢占以平衡不同队列的资源分配<sup>[58]</sup>。但在编写本书时，YARN，Mesos或Kubernetes不支持通用优先级抢占<sup>[60]</sup>。在任务不经常被终止的环境中，MapReduce的设计决策没有多少意义。在下一节中，我们将看看MapReduce的一些替代方法，这些替代方案做出了不同的设计决定。

# 超越MapReduce

尽管MapReduce在20世纪末变得非常流行并被大量炒作，但它只是分布式系统的许多可能的编程模型之一。取决于具体的数据量、数据结构以及处理类型，其他工具可能更适合特定的计算。

尽管如此，我们仍然在讨论MapReduce上花了很多篇幅，因为它是一个有用的学习工具，因为它是分布式文件系统的一个相当清晰和简单的抽象。在这里，简单是指从理解它在做什么的角度来说，而不是容易使用的角度。事实上与容易使用恰恰相反：使用原始的MapReduce API来实现一个复杂的处理任务是相当困难和费力的，例如，你需要从头开始实现全部join算法<sup>[37]</sup>。

针对直接使用的困难，在MapReduce上创建了各种高级编程模型（Pig，Hive，Cascading，Crunch）进一步封装抽象。如果了解MapReduce的工作原理，那么学习这些模型也相当容易，而且它们的高级构造使许多常见的批处理任务更加容易实现。

然而，MapReduce执行模型本身也存在一些问题，这些问题并没有通过增加另一个抽象层次而得到解决，而且在某些类型的处理中表现出糟糕的性能。一方面，MapReduce非常强大：可以使用它来处理那些运行在不可靠的多租户系统上、任务频繁终止的超大规模数据，并且仍然可以完成工作（虽然速度很慢）。另一方面，对于某些类型的处理，其他工具则可能要快几个数量级。

在本章的剩余部分中，我们将看到一些批处理的替代方案。第11章将转向流处理，这可以看作是加速批处理的另一种方法。

## 中间状态实体化

如前所述，每个MapReduce作业都独立于其他任何作业。作业与其他任务的主要联系点是分布式文件系统上的输入和输出目录。如果希望一个作业的输出成为第二个作业的输入，则需要将第二个作业的输入目录配置为与第一个作业的输出目录相同，并且外部工作流调度程序必须仅在第一个作业已经完成后才开始第二个作业。

如果第一个作业的输出是要在组织内部广泛分发的数据集，则此设置是合理的。在这种情况下，需要能够通过名称来引用它，并将其用作多个不同作业（包括由其他团队开发的作业）的输入。将数据发布到分布式文件系统中众所周知的位置可以实现松耦合，这样作业就不需要知道谁在生成输出或者消耗输出（请参见本章前面的“逻辑与布线分离”）。

但是，在很多情况下，我们知道一个作业的输出只能用作另一个作业的输入，这个作业由同一个团队维护。在这种情况下，分布式文件系统上的文件只是中间状态：一种将数据从一个作业传递到下一个作业的方式。在用于构建由50或100个MapReduce作业[29]组成的推荐系统的复杂工作流中，存在很多这样的中间状态。

将这个中间状态写入文件的过程称为实体化（或物化）。我们在第3章“聚合：数据立方体与物化视图”中已经在物化视图的背景下遇到了这个术语，它主要是指提前计算某些操作的结果并将其写入磁盘，而不是在需要时才进行计算。

相比之下，本章开头的日志分析示例使用UNIX管道将一个命令的输出与另一个命令的输入连接起来。管道并不完全实现中间状态，而是只使用一个小的内存缓冲区，逐渐将输出流式传输到输入。

与UNIX管道相比，MapReduce完全实体化中间状态的方法有一些不利之处：

- MapReduce作业只有在前面作业（生成其输入）中的所有任务都完成时才能启动，而通过UNIX管道连接的进程同时启动，输出一旦生成就会被使用。不同机器的差异或不同的负载意味着作业中往往会有一些任务需要比其他任务花费更长的时间。必须等待前面作业中所有任务的完成必然会减慢整个工作流的执行。
- Mapper通常是冗余的：它们只是读取刚刚由reducer写入的同一个文件，并为下一个分区和排序阶段做准备。在许多情况下，mapper代码可能是之前reducer的一部分：如果reducer的输出被分区和排序的方式与mapper输出相同，那么不同阶段的reducer可以直接链接在一起，而不需要与mapper阶段交错。
- 将中间状态存储在分布式文件系统中意味着这些文件被复制到多个节点，对于这样的临时数据来说通常是大材小用了。

## 数据流引擎

为了解决MapReduce的这些问题，开发了用于分布式批处理的新的执行引擎，其中最著名的是Spark<sup>[61, 62]</sup>，Tez<sup>[63, 64]</sup>和Flink<sup>[65, 66]</sup>。它们的设计方式有很多不同之处，但有一个共同点：它们把整个工作流作为一个作业来处理，而不是把它分解成独立的子作业。

由于通过若干个处理阶段明确地建模数据流，所以这些系统被称为数据流引擎。像MapReduce一样，它们通过反复调用用户定义的函数来在单个线程上一次处理一条记录。它们通过对输入进行分区来并行工作，并将一个功能的输出复制到网络上，成为另一个功能的输入。







与MapReduce不同，这些功能不需要严格交替map和reduce的角色，而是以更灵活的方式进行组合。我们称为函数运算符，数据流引擎提供了多种不同的选项来连接一个运算符的输出到另一个的输入：

- 一个选项是通过关键字对记录进行重新分区和排序，就像在MapReduce的shuffle阶段一样（请参阅本章前面的“MapReduce的分布式执行”）。此功能可以像在MapReduce中一样进行排序-合并join和分组。
- 另一个可能性是读取若干个输入，并以相同的方式进行分区，但忽略排序。这节省了分区哈希join的工作，其中记录的分区是重要的，但顺序不相关，因为构建哈希表实现了顺序的随机化。
- 对于广播哈希join，可以将一个运算符的输出发送到join运算符的所有分区。

这种处理引擎的风格源于Dryad<sup>[67]</sup>和Nephele<sup>[68]</sup>这样的研究系统，与MapReduce模型相比，它有几个优点：

- 排序等计算代价昂贵的任务只在实际需要的地方进行，而不是在每个map和reduce阶段之间默认发生。
- 没有不必要的map任务，因为mapper所做的工作通常可以合并到前面的reduce运算符中（mapper不会更改数据集的分区）。
- 由于工作流中的所有join和数据依赖性都是明确声明的，因此调度器知道哪些数据在哪里是必需的，因此它可以进行本地优化。例如，可以尝试将占用某些数据的任务放在与生成它的任务相同的机器上，以便可以通过共享内存缓冲区来交换数据，而不必通过网络复制数据。
- 将运算符之间的中间状态保存在内存中或写入本地磁盘通常就足够了，这比将内容写入HDFS（必须将其复制到多个节点并写入到每个副本所在的磁盘）需要更少的I/O。MapReduce已经将这种优化用于mapper的输出，但是数据流引擎将该思想推广到了所有的中间状态。
- 运算符可以在输入准备就绪后立即开始执行，在下一个开始之前不需要等待前一个阶段全部完成。
- 与MapReduce（为每个任务启动一个新的JVM）相比，现有的Java虚拟机（JVM）进程可以被重用来运行新的运算符，从而减少启动开销。

可以使用数据流引擎来执行与MapReduce工作流相同的计算，并且由于这些优化，通常执行速度会明显加快。由于运算符是map和reduce的一个泛化，相同的处理代码可





以在任一执行引擎上运行：在Pig，Hive或Cascading中所实现的工作流可以通过简单的配置更改从MapReduce切换到Tez或Spark，而无需修改代码<sup>[64]</sup>。

Tez是一个相当轻量的库，它依靠YARN的shuffle服务来实现节点之间的数据复制<sup>[58]</sup>，而Spark和Flink则是包含网络通信层、调度器和面向用户API的大型框架。我们将很快开始讨论这些高级API。

## 容错

将中间状态完全实体化到分布式文件系统的一个优点是持久化，这使得在MapReduce中实现容错变得相当容易：如果一个任务失败了，它可以在另一台机器上重新启动，并从文件系统重新读取相同的输入。

Spark，Flink和Tez避免将中间状态写入HDFS，所以它们采用不同的方法来容忍错误：如果机器发生故障，并且该机器上的中间状态丢失，则利用其他可用的数据重新计算（例如之前的中间阶段，如果可能的话；或原始输入数据，通常在HDFS上）。

为了实现重新计算，框架必须追踪给定数据是如何计算的，使用了哪个输入分区，以及应用了哪个运算符。Spark使用弹性分布式数据集（Resilient Distributed Dataset, RDD）抽象来追踪数据的祖先<sup>[61]</sup>，而Flink对运算符状态建立检查点，从而允许将执行过程中遇到故障的运算符恢复运行<sup>[66]</sup>。

在重新计算数据时，知道计算是否具有确定性非常重要：也就是说，如果给定相同的输入数据，那么运算符是否始终产生相同的输出？如果部分丢失的数据已经发送给下游运算符，这个问题就非常关键。如果运算符重新启动，重新计算的数据与原有的丢失数据不一致，下游运算符会很难解决新旧数据之间的冲突。如果出现非确定性运算符，解决方案通常是下游运算符也终止掉，并在新数据上再次运行它们。

为了避免这种级联故障，最好让运算符具有确定性。但是请注意，非确定性行为很容易意外发生：例如，许多编程语言在迭代哈希表的元素时不能保证任何特定顺序，许多概率和统计算法明确依赖于使用随机数，以及使用系统时钟或外部数据源都是不确定的。为了可靠地从故障中恢复，就需要消除这些不确定性因素，例如通过使用固定的种子产生伪随机数。

通过重新计算数据以实现从故障中恢复并不总是能得到正确的答案：如果中间数据比源数据小得多，或者如果计算量非常大，那么将中间数据转化为文件比重新计算文件的代价要小。





## 关于实体化的讨论

回到UNIX的类比，我们看到MapReduce就像是每个命令的输出写入临时文件，而数据流引擎看起来更像是UNIX管道。尤其Flink是围绕流水线执行的思想而建立起来的，也就是将运算符的输出递增地传递给其他运算符，并且在开始处理之前不等待输入完成。

排序操作不可避免地需要消耗其全部输入，才可以产生输出，因为有可能最后的输入记录具有最小的关键字，因此它必须成为第一个输出记录。任何需要分类的运算符都需要至少暂时地累积状态。但是工作流的许多其他部分可以以流水线方式执行。

当作业完成时，它的输出需要在某个地方实现持久化，以使用户可以找到并使用它，很可能它会再次写入分布式文件系统。因此，在使用数据流引擎时，HDFS上的实体化数据集通常仍是作业的输入和最终输出。和MapReduce一样，输入是不可变的，输出被完全替换。简单总结一下，数据流对MapReduce的改进是，不需要自己将所有中间状态写入文件系统。

## 图与迭代处理

在第2章中的“类图数据模型”中，我们讨论了使用图来建模数据，并使用图查询语言遍历图中的边和顶点。第2章的讨论集中在OLTP的使用方式上：快速查询某些少量符合特定条件的顶点。

在批处理环境中查看图也很有趣，其目标是在整个图上执行某种离线处理或分析。这种需求经常出现在机器学习应用程序（如推荐引擎）或排名系统中。例如，最著名的图分析算法之一是PageRank<sup>[69]</sup>，它试图根据链接至某网页的其他网页来评估该网页的受欢迎程度。它是确定网络搜索引擎结果呈现顺序的标准之一。



像Spark，Flink和Tez这样的数据流引擎（参见本章前面的“中间状态实体化”）通常将运算符作为有向无环图（DAG）排列在作业中。这与图处理并不一样：在数据流引擎中，运算符之间的数据流被构造成一个图，数据本身通常由关系式元组构成；而在图处理中，数据本身具有图的形式。真不幸，又一个由命名引起的混淆！

许多图算法通过一次遍历一个边，将一个顶点与相邻顶点join起来以便传递某种信息，重复该过程直到满足某种条件为止。例如，直到没有更多的边需要遍历，或者直到某些度量值收敛。我们在图2-6中看过一个例子，它通过迭代跟踪那些能够表明某个位置位于哪些其他位置之内（这种算法被称为传递闭包）的边，列出了数据库中所有位于







北美的位置。

可以在分布式文件系统（包含顶点和边的列表文件）中存储图，但是这种“重复直到完成”的想法不能用普通的MapReduce来表示，因为它只执行一次数据传递。这种算法通常需要以迭代方式实现：

1. 外部调度程序运行批处理来执行算法的一个步骤。
2. 当批处理过程完成时，调度器检查遍历是否完成（基于特定完成条件，例如没有更多的边要遍历，或者与最后一次迭代相比的变化低于某个阈值）。
3. 如果尚未完成，则调度程序返回到步骤1并运行另一轮批处理。

这种方法有效，但是用MapReduce来实现却非常低效，主要是因为MapReduce没有考虑算法的迭代性质：即使与上一次迭代相比，只有图的一小部分发生了改变，它也总是读取整个输入数据集并产生一个全新的输出数据集。

## Pregel处理模型

作为对图数据的批处理优化，计算的批量同步并行（*bulk synchronous parallel*, BSP）模型<sup>[70]</sup>已经流行起来，典型系统包括Apache Giraph<sup>[37]</sup>，Spark的GraphX API和Flink的Gelly API<sup>[71]</sup>。由于最早是Google的Pregel论文将这种处理图的方法普及<sup>[72]</sup>，因此它也被称为Pregel模型。

回想一下在MapReduce中，mapper在概念上“发送消息”给reducer调用，因为框架将所有具有相同关键字的mapper输出集中在一起。Pregel中的想法与此类似：一个顶点可以“发送消息”到另一个顶点，通常这些消息沿着图的边被发送。

在每次迭代中，为每个顶点调用函数，将所有发送至该顶点的消息传递给它，就像调用reducer一样。与MapReduce的不同之处在于，Pregel模型的迭代过程中，顶点保存它在内存中的状态，所以函数只需要处理新输入的消息。如果图的某个部分没有收到发送消息，则不需要做任何工作。

这与参与者模型有些相似（请参阅第4章的“分布式Actor model”），可以将每个节点视为参与者，不同之处在于，顶点状态和顶点之间的消息具有容错性和持久性，并且通信以固定的方式进行：每一轮迭代中，框架会将前一次迭代中的所有消息都发送出去。Actor model通常没有这样的时序保证。

## 容错

事实上，顶点只能通过消息传递进行通信（而不是通过彼此间的直接查询）有助于提





高Pregel作业的性能，因为消息可以被批量处理，并且等待通信的次数会减少。唯一的等待是在迭代之间：由于Pregel模型保证在一次迭代中发送的所有消息都会在下一次迭代中被发送，所以先前的迭代必须全部完成，而且所有的消息必须在下一次迭代开始之前复制到网络中。

即使底层网络可能会丢弃、重复或任意延迟消息（请参阅第8章的“不可靠的网络”），但Pregel的实现可以保证在后续迭代中消息在目标顶点只会被处理一次。像MapReduce一样，该框架透明地从故障中恢复，以简化Pregel顶层算法的编程模型。

这种容错方式是通过在迭代结束时定期快照所有顶点的状态来实现的，即将其全部状态写入持久存储。如果某个节点发生故障并且其内存中状态丢失，则最简单的解决方法是将整个图计算回滚到上一个检查点，然后重新开始计算。如果算法是确定性的，并且记录了消息，那么也可以选择性地只恢复丢失的分区（就像我们之前讨论过的数据流引擎）<sup>[72]</sup>。

## 并行执行

顶点不需要知道它运行在哪台物理机器上。当它发送消息到其他顶点时，只需要将消息发送至一个顶点ID。框架对图进行分区，即确定哪个顶点运行在哪个机器上，以及如何通过网络路由消息，以便它们都能到达正确的位置。

由于编程模型一次仅处理一个顶点（有时也称为“像顶点一样思考”），所以框架能够以任意方式划分图。理想情况下，如果顶点之间需要进行大量的通信，那么它们会被分区至同一台机器。但是，找到这样的优化分区是很困难的，通常按照任意分配的顶点ID对图进行分区，而不会尝试将相关的顶点分组在一起。

因此，图算法往往会有很多跨机器通信的开销，中间状态（节点之间发送的消息）往往比原始图大。通过网络发送消息的开销会显著减慢分布式图算法。

出于这样的原因，如果图可以载入到计算机内存中，那么单机（甚至可能是单线程）算法的性能可能要比分布式批处理的性能更好<sup>[73,74]</sup>。即使图大于内存，也可以放在单个计算机的磁盘中，使用GraphChi等框架进行单机处理也是一个可行的选择<sup>[75]</sup>。如果图太大而不适合单个机器，则需要采用像Pregel这样的分布式方法。有效的并行化图算法是一个正在研究中的领域<sup>[76]</sup>。

## 高级API和语言

自MapReduce流行以来，分布式批处理的执行引擎已经逐渐成熟。到目前为止，基础





设施已经足够强大，能够存储和处理超过10000台机器集群中的PB级数据。由于在这种规模下物理操作批处理过程的问题已经或多或少得到了解决，所以问题已经转向其他领域：改进编程模型，提高处理效率，扩大解决的问题域等。

如前所述，由于手工编写MapReduce作业太过耗时费力，因此Hive、Pig、Cascading和Crunch等高级语言和API变得非常流行。随着Tez的出现，这些高级语言还能够移植到新的数据流执行引擎，而无需重写作业代码。Spark和Flink也包含他们自己的高级数据流API，其中很多灵感来自FlumeJava<sup>[34]</sup>。

这些数据流API通常使用关系式构建块来表示计算：将数据集join到某个字段的值上；按关键字对元组进行分组；按条件进行过滤；并通过计数、求和或其他函数来聚合元组。在内部，这些运算使用本章前面讨论过的各种join和分组算法来实现。

除了减少代码的明显优势之外，这些高级接口还允许交互式使用。在这种交互式使用中，可以在shell中逐步编写分析代码，并随时运行以观察计算结果。这种开发方式在探索数据集和测试处理方法时非常有用。这一点很像前面讨论过的UNIX设计哲学。

此外，这些高级接口不仅使提高了系统利用率，而且提高了机器级别的作业执行效率。

## 转向声明式查询语言

与编写执行join的代码相比，指定join作为关系运算符的优点在于，框架可以分析join输入的属性，并自动决定哪个join算法最适合当前的任务。Hive、Spark和Flink利用基于成本的查询优化器来实现这样的功能，甚至还可以改变join顺序，使中间状态数量最小化<sup>[66, 77, 78, 79]</sup>。

join算法的选择会对批处理作业性能产生很大的影响，当然也不需要理解和记住本章讨论过的所有join算法。如果以声明方式指定join，那么可以在应用程序中简单地说明哪些join是必需的，由查询优化器决定如何最佳执行。之前在第2章“数据查询语言”中讨论过这个想法。

但是，在其他方面，MapReduce及后续的数据流引擎与SQL的完全声明式查询模型有很大不同。MapReduce是围绕函数回调的思想构建的：对于每个记录或者一组记录，调用由用户定义的函数（mapper或reducer），并且该函数可以灵活调用任意代码来决定输出。这种方法的优点是，可以利用现有库的大型生态系统来执行数据解析、自然语言分析、图像分析以及运行数值或统计算法等。

轻松运行任意代码是MapReduce之类的批处理系统与MPP数据库的区别所在（参见本







章前面的“对比Hadoop与分布式数据库”）。尽管数据库具有编写用户定义函数的功能，但是使用起来通常会很麻烦，而且与大多数编程语言中广泛使用的程序包管理器和依赖管理系统（例如面向Java的Maven，JavaScript的npm和Ruby的Rubygems）无法很好地集成。

而且，除了join之外，数据流引擎已经证实很多其他情形中，声明性特征也同样具有优势。例如，如果回调函数只包含一个简单的过滤条件，或者只是从一条记录中选择了部分字段，那么在调用每条记录的函数时会有相当多的CPU开销。如果以声明式表示简单的过滤和map操作，那么查询优化器可以利用面向列的存储格式（请参阅第3章的“列式存储”），从磁盘仅读取所需的列。Hive，Spark DataFrames和Impala也使用向量化执行（请参阅第3章的“内存带宽与向量化处理”）：在对CPU高速缓存很友好的内部循环中迭代数据，并避免函数调用。Spark生成JVM字节码<sup>[79]</sup>，Impala使用LLVM为这些内部循环生成本机代码<sup>[41]</sup>。

通过将声明式特征与高级API结合，使查询优化器在执行期间可以利用这些优化方法，批处理框架看起来就更像MPP数据库了（并且能够实现性能相当）。同时，通过具有运行任意代码和读取任意格式数据的可扩展性，它们依然保持了灵活性的优势。

## 不同领域的专业化

尽管能够运行任意代码的可扩展性是有用的，但是标准处理模式不断反复运行的情形也极其常见，因此有必要实现可重用的通用构建模块。传统上，MPP数据库满足了商业智能分析和商业报告的需求，但这只是批处理的诸多应用领域之一。

另一个越来越重要的领域是统计和数值算法，这是诸如分类和推荐系统等机器学习应用所需要的。已经出现了一些可重复使用的实现：例如，Mahout在MapReduce、Spark和Flink之上实现了用于机器学习的各种算法，而MADlib在关系型MPP数据库（Apache HAWQ）中实现了类似的功能<sup>[54]</sup>。

可重用实现对例如k-近邻<sup>[80]</sup>这样的空间算法也是有用的，它在多维空间中搜索与给定项接近的数据项，是一种相似性搜索。近似搜索对于基因组分析算法也很重要，它们需要找到相似但不相同的字符串<sup>[81]</sup>。

批处理引擎正被用于日益广泛的算法领域的分布式执行。随着批处理系统获得更丰富的内置功能和高级声明式运算符，而MPP数据库也变得更具可编程性和灵活性，两者开始变得更加相似：最终，它们都是用于存储和处理数据的系统。





## 小结

本章探讨了批处理这一主题。我们首先查看了诸如awk, grep和sort等UNIX工具, 然后讨论了这些工具的设计理念是如何运用到MapReduce和最新的数据流引擎中。这些设计原则包括: 输入是不可变的, 输出是为了成为另一个(还未知的)程序的输入, 而复杂的问题通过编写“做好一件事”的小工具来解决。

在UNIX世界中, 允许多个程序组合在一起的统一接口是文件和管道; 在MapReduce中, 该接口是分布式文件系统。我们看到数据流引擎添加了自己的管道式数据传输机制, 以避免在分布式文件系统中将中间状态实体化, 而作业的初始输入和最终输出通常仍然是HDFS。

分布式批处理框架需要解决的两个主要问题是:

### 分区

在MapReduce中, mapper根据输入文件块进行分区。mapper的输出被重新分区、排序, 合并成一个可配置数量的reducer分区。这个过程的目的把所有的数据——例如, 具有相同关键字的所有记录都放在同一个地方。

除非必需, 后MapReduce的数据流引擎都尽量避免排序, 但它们采取了大致类似的分区方法。

### 容错

MapReduce需要频繁写入磁盘, 这使得可以从单个失败任务中轻松恢复, 而无需重新启动整个作业, 但在无故障情况下则会减慢执行速度。数据流引擎执行较少的中间状态实体化并保留更多的内存, 这意味着如果节点出现故障, 他们需要重新计算更多的数据。确定性运算符减少了需要重新计算的数据量。

我们讨论了几种MapReduce的join算法, 其中大部分也在MPP数据库和数据流引擎中使用。分区算法的示例包括:

### 排序-合并join

每个将要join的输入都会由一个join关键字的mapper来处理。通过分区、排序与合并, 具有相同关键字的所有记录最终都会进入对reducer的同一个调用。然后这个函数输出join记录。

## 广播哈希join

两个join输入中的某一个数据集很小，所以不需要分区，而完全加载到哈希表中。因此，可以为大数据集的每个分区启动一个mapper，将小数据集的哈希表加载到每个mapper中，然后一次扫描大数据集的一条记录，对每条记录进行哈希表查询。

## 分区哈希join

如果两个join的输入以相同的方式分区（使用相同的关键字，相同的哈希函数和相同数量的分区），则哈希表方法可以独立用于每个分区。

分布式批处理引擎有一个有意限制的编程模型：回调函数（如mapper和reducer）被设定为无状态，并且除了指定输出之外没有外部可见的任何副作用。这个限制使得框架隐藏了抽象背后的一些困难的分布式系统问题，从而在面对崩溃和网络问题时，可以安全地重试任务，并丢弃任何失败任务的输出。如果针对某个分区的多个任务都成功了，则实际上只会有其中一个任务使其输出可见。

得益于这样的框架，在批处理作业中的代码无需考虑容错机制的实现：框架可以保证作业的最终输出与没有发生错误的情况相同（即使实际上各种任务可能不得不重试）。而在线服务在处理用户请求时，将写入数据库作为处理请求的副作用，与之相比，批处理的可靠性语义要强大得多。

批处理作业的显著特点是它读取一些输入数据并产生一些输出数据，而不修改输入。换句话说，输出是从输入派生而来。至关重要的是，输入数据是有界的：数据大小固定已知（例如，它包含一些时间点的日志文件或数据库内容的快照）。因为它是有界的，所以一个作业总是可以知道何时完成了对整个输入的读取，何时作业最终完成。

在下一章中，我们将转向流处理，其中输入是无界的。也就是说，作业的输出是永无止境的数据流。在这种情况下，作业永远不会完成，因为在任何时候都可能有更多的输入进来。我们将看到流处理和批处理在某些方面是相似的，但流数据无界的假设也会深刻改变我们设计系统的方法。

## 参考文献

[1] Jeffrey Dean and Sanjay Ghemawat: “MapReduce: Simplified Data Processing on Large Clusters,” at *6th USENIX Symposium on Operating System Design and Implementation* (OSDI), December 2004.

[2] Joel Spolsky: “The Perils of JavaSchools,” *joelonsoftware.com*, December 25, 2005.



- [3] Shivnath Babu and Herodotos Herodotou: “Massively Parallel Databases and MapReduce Systems,” *Foundations and Trends in Databases*, volume 5, number 1, pages 1-104, November 2013. doi:10.1561/19000000036.
- [4] David J. DeWitt and Michael Stonebraker: “MapReduce: A Major Step Backwards,” originally published at *databasecolumn.vertica.com*, January 17, 2008.
- [5] Henry Robinson: “The Elephant Was a Trojan Horse: On the Death of Map-Reduce at Google,” *the-paper-trail.org*, June 25, 2014.
- [6] “The Hollerith Machine,” United States Census Bureau, *census.gov*.
- [7] “IBM 82, 83, and 84 Sorters Reference Manual,” Edition A24-1034-1, International Business Machines Corporation, July 1962.
- [8] Adam Drake: “Command-Line Tools Can Be 235x Faster than Your Hadoop Cluster,” *aadrake.com*, January 25, 2014.
- [9] “GNU Coreutils 8.23 Documentation,” Free Software Foundation, Inc., 2014.
- [10] Martin Kleppmann: “Kafka, Samza, and the Unix Philosophy of Distributed Data,” *martin.kleppmann.com*, August 5, 2015.
- [11] Doug McIlroy: Internal Bell Labs memo, October 1964. Cited in: Dennis M. Richie: “Advice from Doug McIlroy,” *cm.bell-labs.com*.
- [12] M. D. McIlroy, E. N. Pinson, and B. A. Tague: “UNIX Time-Sharing System:Foreword,” *The Bell System Technical Journal*, volume 57, number 6, pages 1899-1904, July 1978.
- [13] Eric S. Raymond: *The Art of UNIX Programming*. Addison-Wesley, 2003. ISBN:978-0-13-142901-7.
- [14] Ronald Duncan: “Text File Formats — ASCII Delimited Text — Not CSV or TAB Delimited Text,” *ronaldduncan.wordpress.com*, October 31, 2009.
- [15] Alan Kay: “Is ‘Software Engineering’ an Oxymoron?,” *tinlizzie.org*.
- [16] Martin Fowler: “InversionOfControl,” *martinfowler.com*, June 26, 2005.
- [17] Daniel J. Bernstein: “Two File Descriptors for Sockets,” *cr.yp.to*.

- [18] Rob Pike and Dennis M. Ritchie: “The Styx Architecture for Distributed Systems,” *Bell Labs Technical Journal*, volume 4, number 2, pages 146-152, April 1999.
- [19] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung: “The Google File System,” at *19th ACM Symposium on Operating Systems Principles (SOSP)*, October 2003. doi:10.1145/945445.945450.
- [20] Michael Ovsiannikov, Silviu Rus, Damian Reeves, et al.: “The Quantcast File System,” *Proceedings of the VLDB Endowment*, volume 6, number 11, pages 1092-1101, August 2013. doi:10.14778/2536222.2536234.
- [21] “OpenStack Swift 2.6.1 Developer Documentation,” OpenStack Foundation, *docs.openstack.org*, March 2016.
- [22] Zhe Zhang, Andrew Wang, Kai Zheng, et al.: “Introduction to HDFS Erasure Coding in Apache Hadoop,” *blog.cloudera.com*, September 23, 2015.
- [23] Peter Cnudde: “Hadoop Turns 10,” *yahoohadoop.tumblr.com*, February 5, 2016.
- [24] Eric Baldeschwieler: “Thinking About the HDFS vs. Other Storage Technologies,” *hortonworks.com*, July 25, 2012.
- [25] Brendan Gregg: “Manta: Unix Meets Map Reduce,” *dtrace.org*, June 25, 2013.
- [26] Tom White: *Hadoop: The Definitive Guide*, 4th edition. O’Reilly Media, 2015. ISBN: 978-1-491-90163-2.
- [27] Jim N. Gray: “Distributed Computing Economics,” Microsoft Research Tech Report MSR-TR-2003-24, March 2003.
- [28] Márton Trencsényi: “Luigi vs Airflow vs Pinball,” *bytepawn.com*, February 6, 2016.
- [29] Roshan Sumbaly, Jay Kreps, and Sam Shah: “The ‘Big Data’ Ecosystem at LinkedIn,” at *ACM International Conference on Management of Data (SIGMOD)*, July 2013. doi:10.1145/2463676.2463707.
- [30] Alan F. Gates, Olga Natkovich, Shubham Chopra, et al.: “Building a High-Level Dataflow System on Top of Map-Reduce: The Pig Experience,” at *35th International Conference on Very Large Data Bases (VLDB)*, August 2009.
- [31] Ashish Thusoo, Joydeep Sen Sarma, Namit Jain, et al.: “Hive — A Petabyte Scale

Data Warehouse Using Hadoop,” at *26th IEEE International Conference on Data Engineering (ICDE)*, March 2010. doi:10.1109/ICDE.2010.5447738.

[32] “Cascading 3.0 User Guide,” Concurrent, Inc., *docs.cascading.org*, January 2016.

[33] “Apache Crunch User Guide,” Apache Software Foundation, *crunch.apache.org*.

[34] Craig Chambers, Ashish Raniwala, Frances Perry, et al.: “FlumeJava: Easy, Efficient Data-Parallel Pipelines,” at *31st ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, June 2010. doi: 10.1145/1806596.1806638

[35] Jay Kreps: “Why Local State is a Fundamental Primitive in Stream Processing,” *oreilly.com*, July 31, 2014.

[36] Martin Kleppmann: “Rethinking Caching in Web Apps,” *martin.kleppmann.com*, October 1, 2012.

[37] Mark Grover, Ted Malaska, Jonathan Seidman, and Gwen Shapira: *Hadoop Application Architectures*. O’Reilly Media, 2015. ISBN: 978-1-491-90004-8.

[38] Philippe Ajoux, Nathan Bronson, Sanjeev Kumar, et al.: “Challenges to Adopting Stronger Consistency at Scale,” at *15th USENIX Workshop on Hot Topics in Operating Systems (HotOS)*, May 2015.

[39] Sriranjana Manjunath: “Skewed Join,” *wiki.apache.org*, 2009.

[40] David J. DeWitt, Jeffrey F. Naughton, Donovan A. Schneider, and S. Seshadri: “Practical Skew Handling in Parallel Joins,” at *18th International Conference on Very Large Data Bases (VLDB)*, August 1992.

[41] Marcel Kornacker, Alexander Behm, Victor Bittorf, et al.: “Impala: A Modern, Open-Source SQL Engine for Hadoop,” at *7th Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2015.

[42] Matthieu Monsch: “Open-Sourcing PalDB, a Lightweight Companion for Storing Side Data,” *engineering.linkedin.com*, October 26, 2015.

[43] Daniel Peng and Frank Dabek: “Large-Scale Incremental Processing Using Distributed Transactions and Notifications,” at *9th USENIX conference on Operating Systems Design and Implementation (OSDI)*, October 2010.



- [44] “Cloudera Search User Guide,” Cloudera, Inc., September 2015.
- [45] Lili Wu, Sam Shah, Sean Choi, et al.: “The Browsemaps: Collaborative Filtering at LinkedIn,” at *6th Workshop on Recommender Systems and the Social Web (RSWeb)*, October 2014.
- [46] Roshan Sumbaly, Jay Kreps, Lei Gao, et al.: “Serving Large-Scale Batch Computed Data with Project Voldemort,” at *10th USENIX Conference on File and Storage Technologies (FAST)*, February 2012.
- [47] Varun Sharma: “Open-Sourcing Terrapin: A Serving System for Batch Generated Data,” *engineering.pinterest.com*, September 14, 2015.
- [48] Nathan Marz: “ElephantDB,” *slideshare.net*, May 30, 2011.
- [49] Jean-Daniel (JD) Cryans: “How-to: Use HBase Bulk Loading, and Why,” *blog.cloudera.com*, September 27, 2013.
- [50] Nathan Marz: “How to Beat the CAP Theorem,” *nathanmarz.com*, October 13, 2011.
- [51] Molly Bartlett Dishman and Martin Fowler: “Agile Architecture,” at O’Reilly *Software Architecture Conference*, March 2015.
- [52] David J. DeWitt and Jim N. Gray: “Parallel Database Systems: The Future of High Performance Database Systems,” *Communications of the ACM*, volume 35, number 6, pages 85-98, June 1992. doi:10.1145/129888.129894.
- [53] Jay Kreps: “But the multi-tenancy thing is actually really really hard,” *tweetstorm, twitter.com*, October 31, 2014.
- [54] Jeffrey Cohen, Brian Dolan, Mark Dunlap, et al.: “MAD Skills: New Analysis Practices for Big Data,” *Proceedings of the VLDB Endowment*, volume 2, number 2, pages 1481-1492, August 2009. doi:10.14778/1687553.1687576.
- [55] Ignacio Terrizzano, Peter Schwarz, Mary Roth, and John E. Colino: “Data Wrangling: The Challenging Journey from the Wild to the Lake,” at *7th Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2015.
- [56] Paige Roberts: “To Schema on Read or to Schema on Write, That Is the Hadoop

Data Lake Question,” *adaptivesystemsinc.com*, July 2, 2015.

[57] Bobby Johnson and Joseph Adler: “The Sushi Principle: Raw Data Is Better,” at *Strata+Hadoop World*, February 2015.

[58] Vinod Kumar Vavilapalli, Arun C. Murthy, Chris Douglas, et al.: “Apache Hadoop YARN: Yet Another Resource Negotiator,” at *4th ACM Symposium on Cloud Computing (SoCC)*, October 2013. doi:10.1145/2523616.2523633.

[59] Abhishek Verma, Luis Pedrosa, Madhukar Korupolu, et al.: “Large-Scale Cluster Management at Google with Borg,” at *10th European Conference on Computer Systems (EuroSys)*, April 2015. doi:10.1145/2741948.2741964.

[60] Malte Schwarzkopf: “The Evolution of Cluster Scheduler Architectures,” *firmament.io*, March 9, 2016.

[61] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, et al.: “Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing,” at *9th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, April 2012.

[62] Holden Karau, Andy Konwinski, Patrick Wendell, and Matei Zaharia: *Learning Spark*. O’Reilly Media, 2015. ISBN: 978-1-449-35904-1.

[63] Bikas Saha and Hitesh Shah: “Apache Tez: Accelerating Hadoop Query Processing,” at *Hadoop Summit*, June 2014.

[64] Bikas Saha, Hitesh Shah, Siddharth Seth, et al.: “Apache Tez: A Unifying Framework for Modeling and Building Data Processing Applications,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2015. doi:10.1145/2723372.2742790.

[65] Kostas Tzoumas: “Apache Flink: API, Runtime, and Project Roadmap,” *slideshare.net*, January 14, 2015.

[66] Alexander Alexandrov, Rico Bergmann, Stephan Ewen, et al.: “The Stratosphere Platform for Big Data Analytics,” *The VLDB Journal*, volume 23, number 6, pages 939-964, May 2014. doi:10.1007/s00778-014-0357-y.

[67] Michael Isard, Mihai Budiu, Yuan Yu, et al.: “Dryad: Distributed Data-Parallel

Programs from Sequential Building Blocks,” at *European Conference on Computer Systems* (EuroSys), March 2007. doi:10.1145/1272996.1273005.

[68] Daniel Warneke and Odej Kao: “Nephele: Efficient Parallel Data Processing in the Cloud,” at *2nd Workshop on Many-Task Computing on Grids and Supercomputers* (MTAGS), November 2009. doi:10.1145/1646468.1646476.

[69] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd: “The PageRank Citation Ranking: Bringing Order to the Web,” Stanford InfoLab Technical Report 422, 1999.

[70] Leslie G. Valiant: “A Bridging Model for Parallel Computation,” *Communications of the ACM*, volume 33, number 8, pages 103-111, August 1990. doi: 10.1145/79173.79181.

[71] Stephan Ewen, Kostas Tzoumas, Moritz Kaufmann, and Volker Markl: “Spinning Fast Iterative Data Flows,” *Proceedings of the VLDB Endowment*, volume 5, number 11, pages 1268-1279, July 2012. doi:10.14778/2350229.2350245.

[72] Grzegorz Malewicz, Matthew H. Austern, Aart J. C. Bik, et al.: “Pregel: A System for Large-Scale Graph Processing,” at *ACM International Conference on Management of Data* (SIGMOD), June 2010. doi:10.1145/1807167.1807184.

[73] Frank McSherry, Michael Isard, and Derek G. Murray: “Scalability! But at What COST?,” at *15th USENIX Workshop on Hot Topics in Operating Systems* (HotOS), May 2015.

[74] Ionel Gog, Malte Schwarzkopf, Natacha Crooks, et al.: “Musketeer: All for One, One for All in Data Processing Systems,” at *10th European Conference on Computer Systems* (EuroSys), April 2015. doi:10.1145/2741948.2741968.

[75] Aapo Kyrola, Guy Blelloch, and Carlos Guestrin: “GraphChi: Large-Scale Graph Computation on Just a PC,” at *10th USENIX Symposium on Operating Systems Design and Implementation* (OSDI), October 2012.

[76] Andrew Lenharth, Donald Nguyen, and Keshav Pingali: “Parallel Graph Analytics,” *Communications of the ACM*, volume 59, number 5, pages 78-87, May 2016. doi: 10.1145/2901919.





[77] Fabian Hüske: “Peeking into Apache Flink’s Engine Room,” *flink.apache.org*, March 13, 2015.

[78] Mostafa Mokhtar: “Hive 0.14 Cost Based Optimizer (CBO) Technical Overview,” *hortonworks.com*, March 2, 2015.

[79] Michael Armbrust, Reynold S Xin, Cheng Lian, et al.: “Spark SQL: Relational Data Processing in Spark,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2015. doi:10.1145/2723372.2742797.

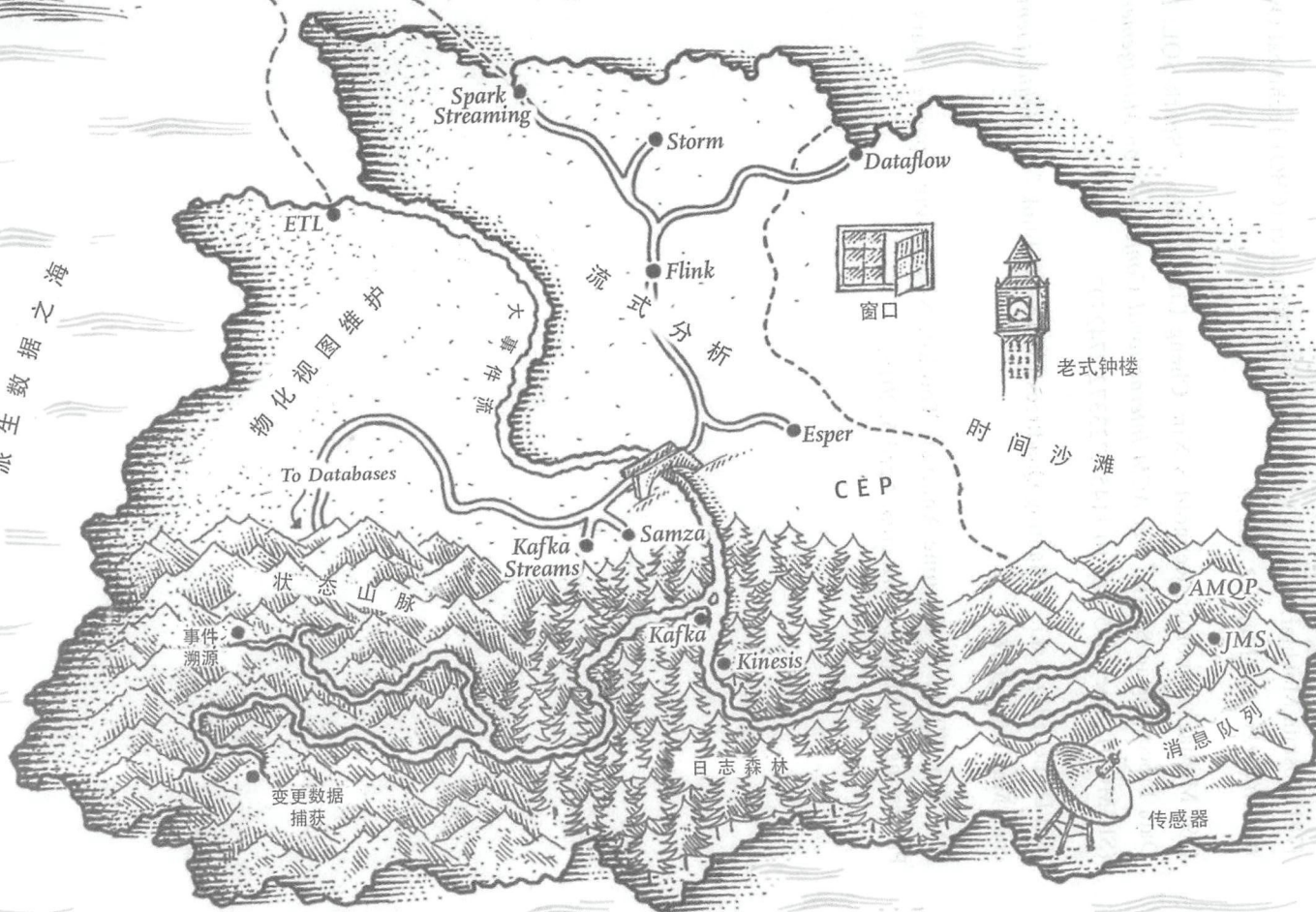
[80] Daniel Blazeovski: “Planting Quadtrees for Apache Flink,” *insightdataengineering.com*, March 25, 2016.

[81] Tom White: “Genome Analysis Toolkit: Now Using Apache Spark for Data Processing,” *blog.cloudera.com*, April 6, 2016.





# 派生数据之海





# 流处理系统

一个可用的复杂系统总是从可用的简单系统演进而来。反过来这话也是正确的：  
从零开始设计的复杂系统从来都用不了，也没办法把它变成可用。

——John Gal, 《系统学》（1975）

第10章讨论了批处理技术，它是读取一组文件作为输入并生成一组新的输出文件的技术。输出是派生数据的一种形式；也就是说，如果需要，可以通过再次运行批处理过程来重新创建数据集。这个看似简单但却强大的想法可以用来创建搜索索引、推荐系统以及分析等。

然而，在第10章中始终存在一个重要的假设：即输入是有界的，是已知的有限大小，所以批处理知道何时读完他们。比如，MapReduce核心的排序操作必须先读取整个输入，然后才能开始生成输出：最后一个输入记录可能是具有最低键的输入记录，因此需要成为第一个输出记录，也正因为如此，不能提前开始输出。

而实际上，有很多数据是无限的，而且随着时间的推移而逐渐到达：用户昨天和今天产生了数据，而明天还将继续产生更多的数据。除非业务停止，否则这个过程永远不会结束，因此数据集永远不会以任何有意义的方式“完成”<sup>[1]</sup>。而批处理器则必须人为地将数据划分为固定时间段的数据块：例如在每天结束时处理一天的数据，或者在每小时结束时处理这一小时的数据。

批处理的问题是，输入的更改只会有一天之后的输出中反映出来，这对于许多没有耐心的用户来说太慢了。为了减少这种延迟，可以更频繁地运行处理，即在每秒钟结束







时（甚至是持续不断地）处理每秒的数据，完全放弃固定的时间片，每当有事件就开始处理。这就是流处理背后的思想。

一般来说，“流”是指随着时间的推移而持续可用的数据。这个概念出现在很多地方：在UNIX的stdin和stdout、编程语言（lazy lists）<sup>[2]</sup>、文件系统API（如Java的FileInputStream）、TCP连接、通过互联网传送音频和视频等。

本章，我们将把事件流视为一种数据管理机制：一种无界的、持续增量处理的方式，对应于上一章所介绍的批处理处理方式。我们将首先讨论如何通过网络来表示、存储和传输流。在“数据库和流”部分，将探讨流和数据库之间的关系。最后，在“处理流”中，将探索连续处理这些流的方法和工具，以及可以用来构建应用程序的方法。

## 发送事件流

在批处理的世界里，作业的输入和输出是文件（或许在分布式文件系统上）。那么等效的流是什么样的呢？

当输入是文件（字节序列）时，第一个处理步骤通常是将其解析为记录序列。在流处理的上下文中，记录通常被称为事件，但它本质上是一回事：一个小的、独立的、不可变的对象，该对象包含某个时间点发生的事情的细节。每个事件通常包含一个时间戳，用于指示事件发生的墙上时间（参见第8章“单调时钟与墙上时钟”）。

例如，发生的事情可能是用户的某些操作，例如浏览页面或下单购买。也可能源于机器设备，例如来自温度传感器的周期性测量或者CPU利用率度量。在第10章的“使用UNIX工具进行批处理”的示例中，Web服务器日志的每一行都是一个事件。

如第4章所述，事件可以被编码为文本字符串或JSON，或者某种二进制形式。通过这种编码方式，可以保存事件，比如将其追加到文件、插入到关系表，或将其写入文档数据库等。它还可以通过网络将事件发送到另一个节点以进行处理。

在批处理中，文件被写入一次，然后可能被多个作业读取。类似地，在流术语中，事件由生产者（也称为发布者或发送者）生成一次，然后可能由多个消费者（订阅者或接收者）<sup>[3]</sup>处理。在文件系统中，文件名标识一组相关记录；在流系统中，相关的事件通常被组合成主题或流。

原则上，通过文件或数据库也可以连接生产者和消费者：生产者将其生成的每个事件写入数据存储，并且每个消费者定期轮询数据存储以检查自上次运行以来出现的事件。这实际上正是批处理在每天结束时处理一天的数据的过程。





但是，如果数据存储不是为这种用途而设计的，那么在延迟时间较短的情况下进行连续处理时，轮询的代价会变得很大。轮询次数越多，返回新事件的请求的百分比越低，开销越高。所以，当新事件出现时，最好通知消费者。

数据库传统上无法很好地支持这种通知机制：关系型数据库通常具有触发器，可以对变化作出反应（例如，将行插入表时），但是他们的功能非常有限，在数据库设计中有点事后考虑的意思<sup>[4,5]</sup>。因此，开发了专门的工具用来提供事件通知。

## 消息系统

向消费者通知新事件的常见方法是使用消息系统：生产者发送包含事件的消息，然后该消息被推送给一个或多个消费者。之前在第4章的“基于消息传递的数据流”只是简单的介绍了这些系统，现在我们来详细讨论这些系统。

像UNIX管道或TCP连接一样，在生产者和消费者之间的直接通信通道将是实现消息传递系统的简单方法。但是，大多数消息传递系统都在这个基本模型上进行扩展。特别是，UNIX管道和TCP仅连接一个发送者与一个接收者，而消息系统允许多个生产者节点将消息发送到同一主题，并允许多个消费者节点接收主题中的消息。

在这种发布/订阅模式中，不同的系统采取了不同的方法，没有一个标准的答案满足所有的目的。为了区分这些系统，提出以下两个问题对区分很有帮助：

1. 如果生产者发送消息的速度比消费者所能处理的快，会发生什么？一般来说，有三种选择：系统丢弃消息；将消息缓存在队列中；或者激活背压，也称为流量控制（即阻止生产者发送更多消息）。例如，UNIX管道和TCP使用背压：他们有一个固定大小的小缓冲区，如果它被填满，发送者将被阻塞，直到接收方将数据从缓冲区中取出（请参阅第8章“网络拥塞与排队”）。

如果消息被缓存在队列中，那么了解队列增长时会发生什么非常重要。如果内存无法容纳所有队列，系统是否会崩溃？还是会将消息写入磁盘？如果是这样，磁盘访问又会如何影响消息传递系统的性能<sup>[6]</sup>？

2. 如果节点崩溃或者暂时离线，是否会有消息丢失？与数据库一样，持久性可能需要写入磁盘和/或结合复制方案（请参阅第7章“复制与持久性”），而这都是有成本的。如果能够接受有时候会丢失消息，那么在同样的硬件上可能获得更高的吞吐量和更低的延迟。

消息丢失是否可以接受在很大程度上取决于应用程序。例如，对于周期性传输的传感器读数 and 指标，偶尔丢失的数据点可能并不重要，因为更新的值马上就会发送过来。





但是，要注意的是，如果大量的消息被丢失，可能相关指标并不会立即发现异常<sup>[7]</sup>。如果你正在对事件进行计数，由于每个丢失的消息都意味着计数器错误，因此它们能够被可靠地传送就显得更重要了。

第10章中所讨论的批处理系统有一个很好的特性是，它们提供了强大的可靠性保证：失败的任务会自动重试，失败任务的部分输出会自动丢弃。这意味着输出与未发生故障时一样，这有助于简化编程模型。本章稍后将研究如何在流上下文中提供类似的保证。

## 生产者与消费者之间的直接消息传递

许多消息系统将生产者直接连接到消费者，而不通过中间节点：

- UDP组播广泛应用于金融行业，例如股票市场等低延迟场景<sup>[8]</sup>。尽管UDP本身是不可靠的，但应用层协议可以恢复丢失的数据包（生产者必须记住它发送的数据包，以便它可以按需要重新发送数据包）。
- 无代理的消息库（如ZeroMQ<sup>[9]</sup>和nanomsg），采取类似的方法，通过TCP或IP多播实现发布/订阅消息传递。
- StatsD<sup>[10]</sup>和Brubeck<sup>[7]</sup>使用不可靠的UDP消息传递来收集网络中所有机器的指标并对其进行监控（在StatsD协议中，只有接收到所有消息时，计数器指标才是正确的；使用UDP只能使计数近似准确<sup>[11]</sup>。另请参阅第8章“TCP与UDP”）。
- 如果消费者在网络上公开服务，则生产者可以直接发出HTTP或RPC请求（请参阅第4章“基于服务的数据流：REST和RPC”）以将消息推送给消费者。这正是webhooks背后的想法<sup>[12]</sup>：一个服务的回调URL被注册到另一个服务中，并且每当事件发生时都会向该URL发出请求。

在其设计的目标场景下，这些直接消息传递方式运行效果不错，但是它们通常都要求应用程序代码意识消息丢失的可能性。它们只能支持有限的容错：即使协议可以检测并重新传输在网络中丢失的数据包，但通常还是假定生产者和消费者需要一直在线。

如果消费者处于离线状态，则可能会遗漏当他们掉线时发送的消息。有些协议允许生产者重试失败的消息传递，然而，如果生产者崩溃，则这种方法可能会失败，从而丢失了本应该重试的消息缓冲区。

## 消息代理

一种广泛使用的替代方法是通过消息代理（也称为消息队列）发送消息，消息代理实质上是一种针对处理消息流而优化的数据库<sup>[13]</sup>。它作为服务器运行，生产者和消费者作为客户端连接到它。生产者将消息写入代理，消费者通过从消息代理那里读取消息







来接收消息。

通过将数据集中在代理中，这些系统可以更容易地适应不断变化的客户端（连接、断开连接和崩溃），而持久性问题则被转移到代理那里。一些消息代理只将消息保存在内存中，而另一些消息代理（取决于配置）将其写入磁盘，以便在代理崩溃的情况下不会丢失消息。对于速度慢的消费者，他们通常允许无限队列（而不是丢弃消息或背压），不过这种选择也可能取决于配置。

排队的结果也通常导致消费者以异步方式工作：生产者发送消息时，它通常只等待代理确认它已经缓存了消息，而不会等待消息被消费者处理。向消费者的交付发生在将来某个不确定的时间点——通常是在几分之一秒内，但如果存在队列积压，有时会有很明显的延迟。

## 消息代理与数据库对比

一些消息代理甚至可以使用XA或JTA参与两阶段提交协议（请参阅第9章的“实践中的分布式事务”）。这个特性使它们在本质上与数据库非常相似，虽然消息代理和数据库之间仍然存在着重要的实际差异：

- 数据库通常会保留数据直到被明确要求删除，而大多数消息代理在消息成功传递给消费者时就自动删除消息。这样的消息代理不适合长期的数据存储。
- 由于消息代理很快删除了消息，多数消息系统会假定当前工作集相当小，即队列很短。如果因为消费者速度很慢，而使代理需要缓存很多消息的话（如果内存无法容纳所有的消息，可能会将部分消息唤出到磁盘），那么每个消息就需要更长的时间来处理，整个吞吐量可能会因此降低<sup>[6]</sup>。
- 数据库通常支持二级索引和各种搜索数据的方式，而消息代理通常支持某种方式订阅匹配特定模式的主题。这些机制虽然是不同的，但本质上都是让客户端可以选择它们想要了解的部分数据。
- 查询数据库时，结果通常基于数据的时间点快照。如果另一个客户端随后向数据库写入更改查询结果的内容，那么第一个客户端不会发现之前的结果已经过期（除非它重复查询或轮询更改）。相比之下，消息代理不支持任意的查询，但是当数据发生变化时（即新消息可用时），它们会通知客户端。

这是消息代理的传统观点，体现在像JMS<sup>[14]</sup>和AMQP<sup>[15]</sup>这样的标准中，并有很多系统实现，包括RabbitMQ, ActiveMQ, HornetQ, Qpid, TIBCO Enterprise Message Service, IBM MQ, Azure Service Bus, 以及Google Cloud Pub / Sub等<sup>[16]</sup>。





## 多个消费者

当多个消费者读取同一个主题中的消息时，有两种主要的消息传递模式，如图11-1所示。

### 负载均衡式

每一条消息都只被传递给其中一个消费者，所以消费者可以共享主题中处理消息的工作。代理可以任意分配消息给消费者。当处理消息的代价很高时，此模式非常有用，因此希望能够添加消费者来并行处理消息（在AMQP中，可以通过让多个客户端使用同一个队列消费来实现负载均衡，而在JMS中，它称为共享订阅）。

### 扇出式

每条消息都被传递给所有的消费者。扇出允许几个独立的消费者各自“收听”相同的消息广播，而不会相互影响，流相当于多个读取相同输入文件的不同批处理作业（此功能由JMS中的主题订阅提供和AMQP中的交换绑定）。

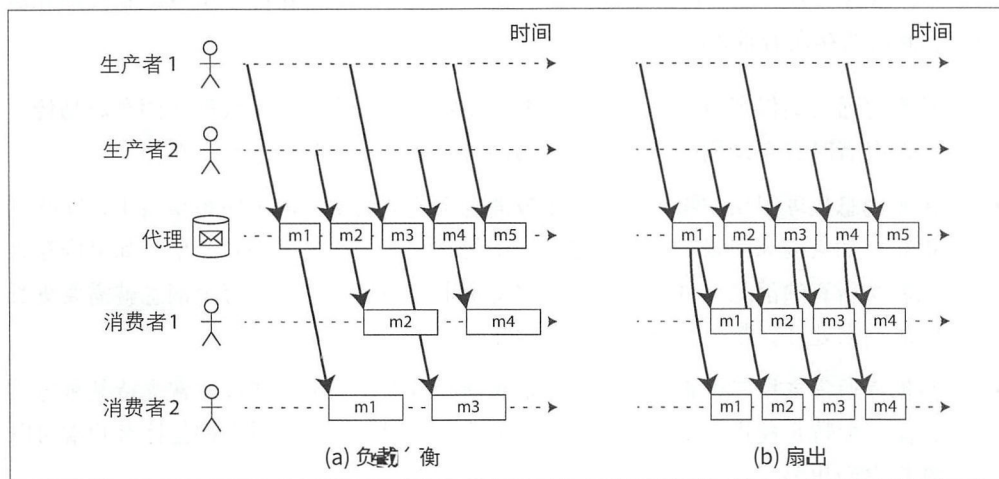


图11-1：(a)负载均衡：多个消费者共同分担消费同一个主题；(b)扇出：将消息传递给多个消费者

这两种模式可以组合使用。例如，两个独立的消费者群组可以各自订阅一个主题，使得每个组都能共同接收所有消息，但是在每个组内，每一条信息只有一个节点接收。

## 确认和重新传递

消费者可能会随时崩溃，所以可能会发生下面这些情况：代理向消费者传递消息，但消费者从不处理消息，或者在崩溃之前只对消息进行了部分处理。为了确保消息不会



丢失，消息代理使用确认：客户端必须在处理完消息后显式地告诉代理，以便代理可以将其从队列中移除。

如果与客户端的连接关闭或超时，而代理没有收到确认，则认为消息未处理，因此它将消息重新传递给另一个消费者（请注意，消息可能实际上已经完全处理，但是确认消息在网络传输过程中丢失。处理这种情况需要原子提交协议，正如在第9章中所讨论的“实践中的分布式事务”那样）。

当与负载均衡结合时，这种重新传递行为对消息的排序会产生一个有趣的影响。在图11-2中，消费者通常按照生产者发送的顺序处理消息。然而，消费者2在处理消息m3时崩溃，与此同时消费者1正在处理消息m4。未确认的消息m3随后被重新发送给消费者1，结果消费者1按照m4，m3，m5的顺序处理消息。因此，m3和m4不是以它们被生产者1发送顺序传递的。

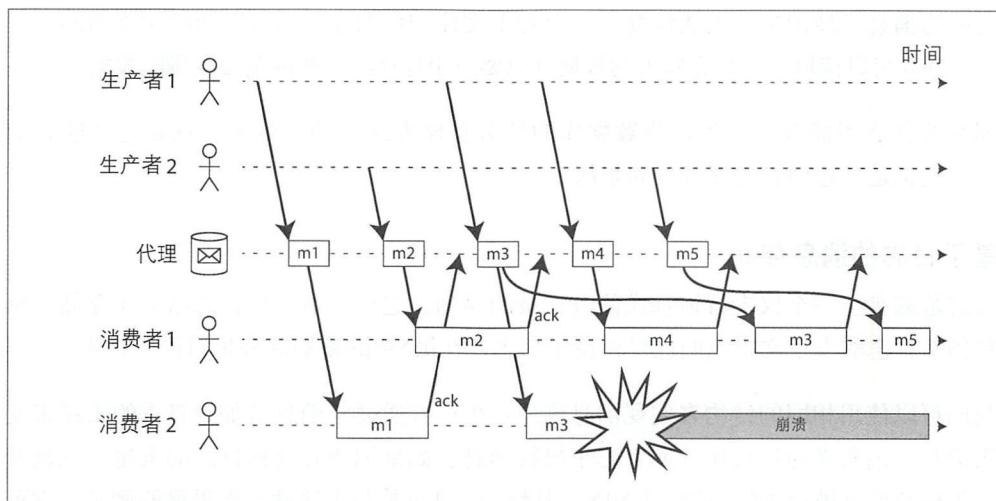


图11-2：消费者2在处理m3时崩溃，因此稍后会再次向消费者1发送m3

即使消息代理试图保留消息的顺序（如JMS和AMQP标准所要求的），负载均衡与重新传递的组合也不可避免地导致消息被重新排序。为了避免此问题，可以为每个消费者使用单独的队列（即不使用负载均衡功能）。如果消息彼此完全独立，消息重新排序就不成问题，但是如果消息之间存在因果依赖关系，那么它就成为很重要的问题了，我们将在本章后面继续讨论。

## 分区日志

通过网络发送数据包或者向网络服务发送请求通常都是瞬间的操作，不会留下永久的



痕迹。尽管可以永久记录（使用数据包捕获和日志记录），但通常不会这么去做。即使是将消息持久地写入磁盘的消息代理，在将消息传递给消费者之后，也会很快将其删除，因为消息代理是基于瞬间的消息传递思维构建的。

数据库和文件系统采取相反的方式：在有人明确选择删除它之前，任何写入数据库或文件的内容通常都期望是永久保存。

思维方式上的这种差异对于如何创建派生数据有很大的影响。正如在第10章讨论的一样，批处理过程的一个关键特征是，可以反复地运行它们，尝试处理步骤，而且没有损坏输入的风险（因为输入是只读的）。但是AMQP / JMS风格的消息系统则不是这样的：如果确认后从代理中删除了消息，就无法再次接收消息，因此不能再次运行同一个消费者并期望得到相同的结果。

如果将新的消费者添加到消息系统，通常它只会开始接收在它注册后发送的消息；任何之前的消息已经消失，无法恢复了。而对于文件和数据库，则可以随时添加新的客户端，并且可以读取以前任意写入的数据（只要应用程序没有明确覆盖或删除数据）。

那么为什么不能混合使用，将数据库的持久存储方法与消息传递的低延迟功能相结合？这正是日志消息代理背后的想法。

## 基于日志的消息存储

日志是磁盘上一个仅支持追加式修改记录的序列。之前在第三章日志结构化存储引擎和预写日志的上下文中我们已经讨论了日志，在第5章的复制部分也讨论了日志。

我们可以使用相同的结构来实现消息代理：生产者通过将消息追加到日志的末尾来发送消息，消费者通过依次读取日志来接收消息。如果消费者读到日志的末尾，它就开始等待新消息被追加的通知。UNIX工具`tail -f`正是基于这种工作思路的例子，它可以监视修改文件的尾部。

为了突破单个磁盘所能提供的带宽吞吐的上限，可以对日志进行分区（参阅第6章）。不同的节点负责不同的分区，使每个分区成为一个单独的日志，并且可以独立于其他分区读取和写入。然后将主题定义为一组分区，他们都携带相同类型的消息。这种方法如图11-3所示。

在每个分区中，代理为每个消息分配一个单调递增的序列号或偏移量（在图11-3中，框中的数字是消息偏移量）。这样的序列号是非常有意义，因为分区只能追加，所以分区内的消息是完全有序的。不同分区之间则没有顺序保证。

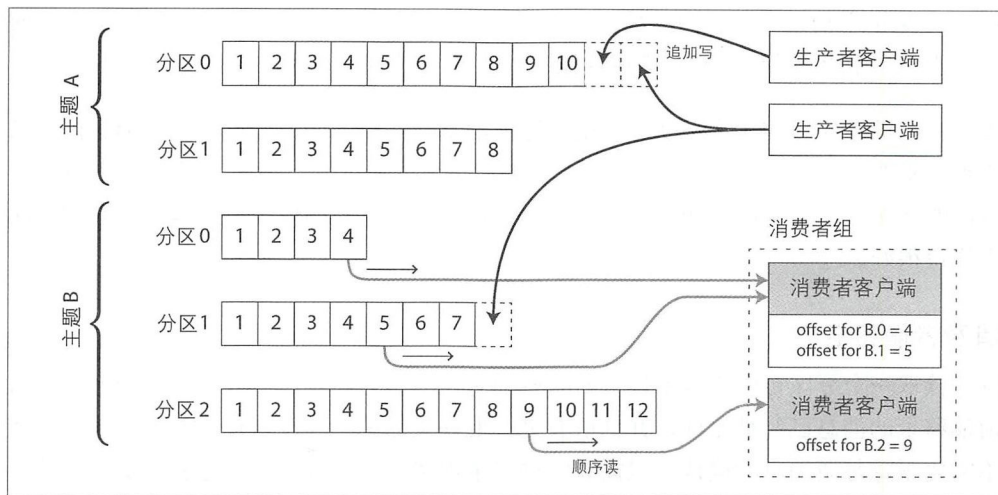


图11-3：生产者通过将消息追加到基于主题-分区的文件，消费者依次读取这些文件

Apache Kafka<sup>[17,18]</sup>，Amazon Kinesis Streams<sup>[19]</sup>和Twitter DistributedLog<sup>[20,21]</sup>都是这种方式工作的基于日志的消息代理系统。Google Cloud Pub/Sub在架构上相似，但是开放JMS风格的API而不是日志抽象<sup>[16]</sup>。尽管这些消息代理将所有消息写入磁盘，但通过在一组消费者之间进行分区，能够实现每秒数百万条消息的吞吐量，并且通过复制消息实现了容错性<sup>[22,23]</sup>。

## 对比日志与传统消息系统

因为多个消费者可以独立地读取日志而不会相互影响，读取消息不会将其从日志中删除，因此基于日志的方法很自然地支持扇出式消息传递。为了在一组消费者之间实现负载均衡，代理可以将整个分区分配给消费者组中的节点，而不是将单个消息分配给消费者客户端。

每个客户端都会使用分配给它所在分区中的所有消息。通常，当消费者被分配了一个日志分区时，它将以直接的单线程方式顺序读取分区中的消息。这种粗粒度的负载均衡方法有一些缺陷：

- 因为同一分区内的消息将被传递到同一节点<sup>注1</sup>，所以消费一个主题的节点数最多等于该主题中的日志分区数。

注1：可以创建一个负载均衡方案，在这种方案中，两个消费者通过读取全部消息来共享处理分区的工作，但其中一个只考虑偶数偏移量的消息，而另一个处理奇数偏移量的消息。或者，可以将消息处理扩展到线程池，但这种方法会使消费者偏移管理变得复杂。通常，单线程处理分区是优选的，并且可以通过使用更多的分区来增加并行性。

- 如果单个消息处理缓慢，则会阻碍该分区中的后续消息的处理（一种队头阻塞的形式，请参阅第1章的“描述性能”）。

因此，在消息处理的代价很高，希望在逐个消息的基础上并行处理，而且消息排序又不那么重要的情况下，JMS / AMQP类型的消息代理更可取。另一方面，在消息吞吐量高的情况下，每个消息处理速度快，消息顺序又很重要的情况下，基于日志的方法工作得很好。

## 消费者偏移量

顺序读取一个分区可以很容易地判断哪些消息已经被处理：所有偏移量小于消费者当前偏移量的消息已经被处理，并且所有更大偏移量的消息还没有被看到。因此，代理不需要跟踪每条消息的确认，只需要定期记录消费者的偏移量。在这种方法中，减少的记录开销以及可以使用批处理和流水线操作的机会将有助于提高基于日志的系统的吞吐量。

实际上，此偏移量与主从复制数据库常见的日志序列号非常相似，在第5章的“配置新的从节点”中讨论了这种情况。在数据库复制中，日志序列号允许从节点断开连接后，重新连接到主节点，并在不跳过任何写入的情况下恢复复制。这里使用了完全相同的原理：消息代理的行为就像一个主节点数据库，消费者就像一个从节点。

如果消费者节点失败，则消费者组中的另一个节点将被分配到失败的消费者分区，并以最后记录的偏移量开始使用消息。如果消费者已经处理了后续的消息，但还没有记录它们的偏移量，那么在重新启动后这些消息将被再次处理。本章后面将讨论处理这个问题的方法。

## 磁盘空间使用

如果持续不断地追加日志，磁盘空间最终将被耗尽。为了回收磁盘空间，日志实际上是被分割成段，并且不时地将旧段删除或归档保存。稍后我们将讨论一种更复杂的释放磁盘空间的方法。

这就意味着，如果一个消费者的速度慢到难以跟上消息产生的速度，并且远远落后以至于消费者偏移量指向了已经被删除的片段，那么消费者将会错过一些消息。实际上，日志实现了一个有限大小的缓冲区，当缓冲区变满时，旧的消息就被丢弃，该缓冲区也被称为循环缓冲区或环形缓冲区。由于该缓冲区在磁盘上，因此它可以非常大。

让我们来做一下粗略的计算。在撰写本文时，一个典型的大容量硬盘为6TB，顺序写



入吞吐量150MB/s。如果以最快速度写入消息，则大约11小时就可以填满磁盘。因此，磁盘可以缓存11小时的消息，之后它将开始覆盖旧的消息。即使使用多个硬盘和多台机器，这个比率也是一样的。实际的部署中，很少会达到磁盘的满写带宽，所以日志通常可以保存几天甚至几周的消息缓冲区。

不管保留多长时间的消息，因为每个消息都被写入到磁盘，因此日志的吞吐量基本保持不变<sup>[18]</sup>。这种行为与将消息默认保存在内存中，仅当队列变得过大时才将它们写入磁盘的消息传递系统相比，差异明显：当队列很短的时候这些系统是很快的，当开始写入磁盘时，会变得很慢，因此吞吐量取决于保留的历史记录数量。

## 当消费者跟不上生产者时

本章前面“消息传递系统”的开始部分，讨论了消费者无法跟上生产者发送消息的速度时三种选择：丢弃消息，缓冲或应用背压。在这个分类法中，基于日志的方法是一种缓冲形式，它具有较大但固定大小的缓冲区（受可用磁盘空间的限制）。

如果消费者落后得太多，以至于其所需的信息比保留在磁盘上的信息还要旧，那么它将无法读取这些信息，所以代理有效地丢弃缓冲区容量不能容纳的旧消息。可以监控消费者落后日志头部的距离，并在落后明显时会发出警报。由于缓冲区很大，因此有足够的时间让操作员修复缓慢的消费者，并允许它在开始丢失消息前赶上。

即使消费者确实落后太多，并且开始丢失信息，也只有该消费者受到影响；它不会中断其他消费者的服务。这是一个非常大的运营优势：可以通过实验性的方式使用生产日志进行开发、测试或调试，而不必担心中断生产服务。当消费者关闭或崩溃时，它会停止消耗资源，唯一留下的就是消费者偏移量。

这种行为也与传统的消息代理不同，在传统的消息代理中，需要小心地删除消费者已经关闭的任何队列，否则他们将继续不必要地积累消息，并占用其他活动消费者的内存。

## 重新处理信息

我们之前提到过，使用AMQP和JMS风格的消息代理时，由于会导致消息在代理上被删除，因此处理和确认操作可视为带有一定的破坏性。另一方面，在基于日志的消息代理中，使用消息更像是从文件读取：这是只读操作，并不会更改日志。

除了消费者的任何输出之外，处理的唯一副作用是消费者偏移量前移了。但是偏移量在消费者的控制之下，因此在必要时可以轻松地进行操作。例如，可以用昨天的

偏移量启动一个消费者的副本，并将输出写到不同的位置，以便重新处理最后一天的消息。可以通过改变处理代码多次重复此操作。

这个特点使得基于日志的消息系统更像上一章的批处理过程，其中派生数据通过可重复的转换过程与输入数据明确分离。它支持更多的实验性尝试，也更容易从错误和故障中进行恢复，从而成为集成数据流的不错选择<sup>[24]</sup>。

## 数据库与流

我们已经在消息代理和数据库之间进行了一些比较。尽管传统上他们被认为是单独的工具类别，但是我们看到基于日志的消息代理已经成功地从数据库中得到启发，并将其应用于消息传递。也可以反过来：从消息传递和流中获取一些启发，并将他们应用到数据库。

之前曾经说过，事件是某个时刻发生的事情的记录。发生的事情可能是用户操作（例如输入搜索查询）或传感器读取，但也可能是写入数据库。将内容写入数据库的事实是一个可以被捕获、存储和处理的事件。这一观察结果表明，数据库和数据流之间的联系比磁盘上日志的物理存储更紧密，这点非常重要。

实际上，复制日志（请参阅第5章的“复制日志的实现”）是数据库写入事件的流，由主节点在处理事务时生成。从节点将写入流应用于他们自己的数据库副本，从而最终得到相同数据的准确副本。复制日志中的事件描述了数据变化。

我们还在第9章的“全序广播”中讨论过状态机复制原理，该原理指出：如果每个事件代表对数据库的写入，并且每个副本按相同的顺序处理相同的事件，则所有副本最后都将收敛于相同的最终状态（处理事件被假设为确定性的操作）。这只是事件流的另一种情况！

在本节中，我们将首先看看异构数据系统中出现的问题，然后探讨如何通过将事件流的想法引入到数据库来解决这个问题。

## 保持系统同步

正如在本书中所看到的，没有一个系统能够满足所有的数据存储、查询和处理需求。在实践中，大多数重要的应用程序都需要结合多种不同的技术来满足需求：例如，使用OLTP数据库来为用户请求提供服务，使用缓存来加速常见请求，使用全文索引处理搜索查询，以及使用数据仓库用于分析。每一个技术都有自己的数据副本，以自己的表示方法存储，并且针对自己的设计目标而优化。



由于相同或相关的数据出现在多个不同的地方，因此他们需要保持相互同步：如果数据库中的某个项更新，则也需要在缓存、搜索索引和数据仓库中进行更新。对于数据仓库，这种同步通常由ETL进程执行（请参阅第3章的“数据仓库”），一般通过获取数据库的完整副本，对其进行转换并将其批量加载到数据仓库中——换句话说，就是批处理。同样，在第10章的“批处理工作流的输出”中介绍了如使用批处理过程创建搜索索引、推荐系统和其他派生数据系统。

如果定期的完整数据库转储过于缓慢，有时使用的替代方法是双重写入，其中程序代码在数据更改时显式地写入每个系统。例如，首先写入数据库，然后更新搜索索引，然后使缓存条目失效（或者甚至同时执行这些写入）。

但是，双重写入有一些严重的问题，其中一个就是图11-4所示的竞争条件。在这个例子中，两个客户同时想要更新项X：客户端1想要将值设置为A，客户端2想要将其设置为B。两个客户端首先将新值写入数据库，然后将其写入搜索索引。由于时机不凑巧，这些请求交叉了：数据库首先看到来自客户端1的写入，将值设置为A，然后看到来自客户端2的写入，将值设置为B，因此数据库中的最终值为B。而搜索索引首先看到来自客户端2的写入，然后才是客户端1的写入，所以搜索索引中的最终值是A。这两个系统将永远不一致，即使目前还没有发生错误。

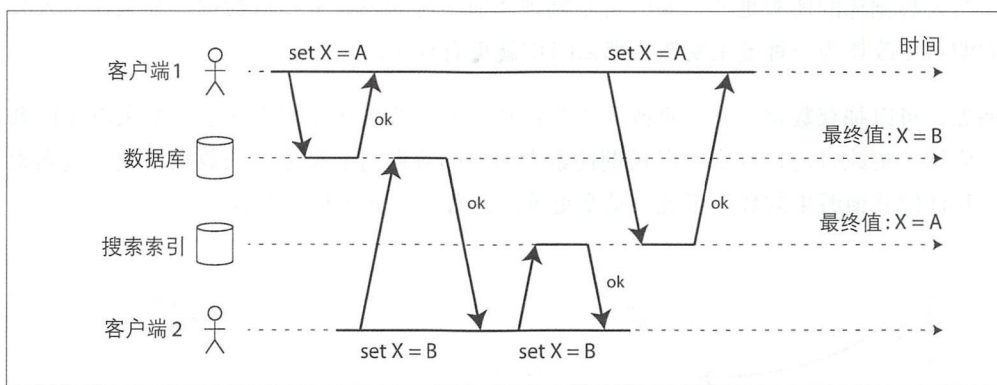


图11-4：数据库中X首先被设置为A，然后被设置为B，而搜索索引则以相反的顺序到达

除非有一些额外的并发检测机制，例如在第5章的“检测并发写”中讨论的版本向量，否则甚至不会注意到发生了并发写入，一个值悄悄地覆盖另一个值。

双重写入的另一个问题是其中一个写入可能会失败，而另一个却成功了。这其实是一个容错问题而不是并发问题，但也会造成两个系统相互不一致的结果。确保他们都成功或者都失败属于原子提交范畴，这个问题解决起来代价很大（请参阅第9章“原子提交与两阶段提交”）。



如果只有一个复制的数据库与一个主节点，那么主节点确定写入的顺序，所以状态机复制方法适用于数据库的副本。然而，在图11-4中没有一个单独的主节点：数据库可能有一个主节点，而搜索索引可能也有其主节点，但是两者都不会跟随对方的主节点，所以冲突可能就发生了（参阅第5章“多主节点复制”）。

如果实际上只有一个主节点（例如数据库），并且如果可以使搜索索引成为数据库的从节点，情况会更好。但这在实践中可能吗？

## 变更数据捕获

大多数数据库的复制日志的问题在于，它们长期以来被认为是数据库的内部实现细节，而不是公开的API。客户端应该通过其数据模型和查询语言来查询数据库，而不是分析复制日志并尝试从中提取数据。

数十年来，许多数据库根本没有任何一种详细的方法来获取所写入的变更日志。由于这个原因，很难将数据库中所做的所有更改复制到不同的存储技术，如搜索索引，缓存或数据仓库。

最近，人们对变更数据捕获（Change Data Capture，CDC）越来越感兴趣。CDC记录了写入数据库的所有更改，并以可复制到其他系统的形式来提取数据。如果在写入时立即将更改作为一种流来发布，那么CDC就更有意思了。

例如，可以捕获数据库中的更改并不断将相同的更改应用于搜索索引。如果以相同顺序应用于更改日志，那么可以预期搜索索引中的数据与数据库中的数据匹配。搜索索引和任何其他派生的数据系统只是变更流的消费者，如图11-5所示。

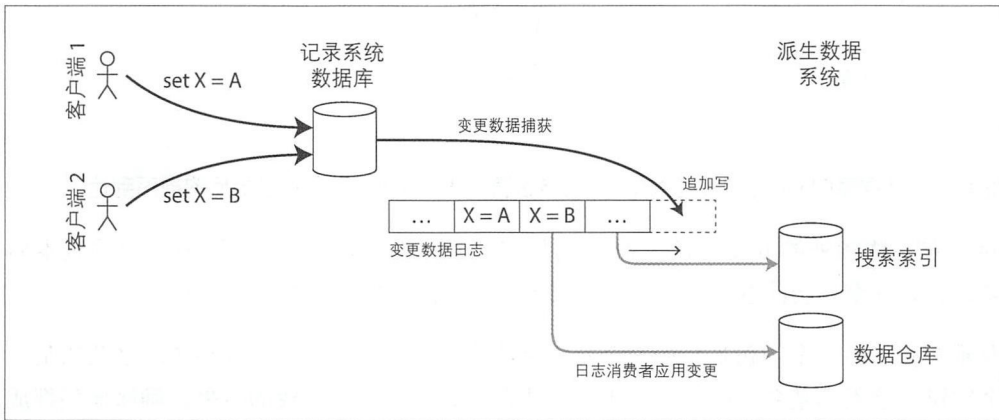


图11-5：将数据按顺序写入数据库，然后按照相同的顺序将这些变化应用到其他系统



## 实现变更数据捕获

我们可以调用日志消费者的派生数据，如第三部分介绍中所述：存储在搜索索引和数据仓库中的数据只是记录系统中数据的另一个视图。变更数据捕获机制可以确保对记录系统所做的所有更改都反映在派生数据系统中，以便派生系统具有数据的准确副本。

从本质上讲，变更数据捕获使得一个数据库成为主节点（从中捕获变化的数据库），并将其他变成从节点。由于基于日志的消息代理保留了消息的排序，因此它非常适合从源数据库传输更改事件（避免了图11-2的重新排序问题）。

数据库触发器可以通过注册触发器来实现变更数据捕获（请参阅第5章“基于触发器的复制”），这些触发器可观测到数据表的所有更改，并将相应的条目添加到更改日志表中。但是，他们往往不太稳定，并且有着非常大的性能开销。解析复制日志是一种更健壮的方法，但它也带来了挑战，比如处理模式更改。

基于该想法的LinkedIn Databus<sup>[25]</sup>，Facebook Wormhole<sup>[26]</sup>和Yahoo! Sherpa<sup>[27]</sup>已在大规模环境下得到部署。Bottled Water使用解码预写日志的API来实现PostgreSQL的CDC<sup>[28]</sup>，Maxwell和Debezium通过解析binlog<sup>[29-31]</sup>为MySQL数据库做类似的事情，Mongoriver读取MongoDB的oplog<sup>[32,33]</sup>，Oracle GoldenGate也提供类似的功能<sup>[34,35]</sup>。

像消息代理一样，变更数据捕获通常是异步的：记录数据库系统不会在提交更改之前等待应用于消费者。这种设计具有的操作上的优势是，添加缓慢的消费者不会对记录系统造成太大影响，但是它的缺点是所有复制滞后导致的问题在这里全部适用（请参阅第5章“复制滞后问题”）。

## 初始快照

如果有了数据库所有更改的日志，就可以通过replay日志来重建数据库的整个状态。然而，在许多情况下，永久保留所有更改将会需要太多的磁盘空间，并且replay将花费太长时间，因此日志需要被截断。

例如，构建新的全文索引需要整个数据库的完整副本，仅仅应用最近更改的日志还不够，因为它会丢失最近未更新的项目。因此，如果没有完整的日志历史记录，则需要从一致的快照开始，正如之前在第5章“配置新的从节点”所讨论的。

数据库的快照必须与更改日志中的已知位置或偏移量相对应，以便在快照处理完成后，知道在哪一点开始应用更改。一些CDC工具集成了此快照功能，而另外一些则需要手动操作。





## 日志压缩

如果只能保留有限的日志历史记录，则每次需要添加新的派生数据系统时都需要执行快照过程。但是，日志压缩提供了一个很好的选择。

之前在第3章的“哈希索引”部分的日志结构存储引擎（见图3-2中的例子）中，讨论了日志压缩功能。原理很简单：存储引擎定期查找具有相同key的日志记录，丢弃所有的重复项，并且只保留每个key的最新的更新。这个压缩和合并的过程是在后台运行的。

在日志结构存储引擎中，具有特殊空值（逻辑删除）的更新实际意味着key被删除了，从而导致它在日志压缩过程中被清除。但只要key不被覆盖或删除，它就永远留在日志中。这种压缩日志所需的磁盘空间仅取决于数据库的当前内容，而不取决于数据库中发生的写入次数。如果相同的key经常被覆盖，则以前的值最终将被垃圾回收，只有最新的值将被保留。

在基于日志的消息代理和变更数据捕获上下文中，相同的想法也适用。如果CDC系统设置为每个更改都有一个key，并且每个key的更新都替换了之前的值，那么对一个特定的key，仅保留最近一次写入就足够了。

现在，无论何时重建派生数据系统（如搜索索引），都可以从日志压缩主题的偏移量0的位置启动一个新的消费者，然后依次扫描日志中的所有消息。日志确保包含数据库中每个key的最新值（也可能是一些较旧的值）。换句话说，可以使用它来获取数据库内容的完整副本，而无需生成CDC源数据库的另一个快照。

Apache Kafka支持此日志压缩功能。正如将在本章后面看到的，这样消息代理可用于永久存储，而不仅仅是用于临时消息传递。

## 对变更流的API支持

越来越多的数据库开始支持将变更流作为标准接口，而不是那些典型的改进和反向工程的CDC努力。例如，RethinkDB支持订阅查询结果发生变化的通知，Firebase<sup>[37]</sup>和CouchDB<sup>[38]</sup>的数据同步基于change feed并同时提供给应用层，而Meteor使用MongoDB oplog来订阅数据更改消息并更新用户界面<sup>[39]</sup>。

VoltDB支持事务以流的形式连续地从数据库中导出数据<sup>[40]</sup>。数据库将关系数据模型中的输出流表示为表，该表支持事务插入元组，但不支持查询。输出流包含了向该特殊表提交写事务的元组日志，并严格按照事务提交顺序排序。外部消费者可以异步使用此日志并使用它来更新派生数据系统。







Kafka Connect<sup>[41]</sup>致力于将广泛的数据库系统变更数据采集工具与Kafka集成。一旦更改事件流汲取到Kafka中，它就可以用来更新派生数据系统，比如搜索索引，也可以用于本章稍后讨论的流处理系统。

## 事件溯源

这里所讨论的想法与事件溯源之间有一些相似之处，它是一种在领域驱动设计（DDD）社区中开发的技术。接下来，我们将简要讨论事件溯源，它包含了一些和流系统相关并且有用的想法。

与变更数据捕获类似，事件溯源涉及到将所有对应用程序状态的更改保存为更改事件的日志。最大的区别在于事件溯源在不同抽象层次上应用了这个想法：

- 在变更数据捕获中，应用程序以数据可变方式来操纵数据库，例如自由地更新和删除记录。从数据库中提取较低级别的变更日志（例如，通过解析复制日志），从而确保从数据库提取的写入顺序与实际写入的顺序相匹配，从而避免图11-4中的竞争条件。写入数据库的应用程序不需要知道CDC正在发生。
- 在事件溯源中，应用程序逻辑是基于写入事件日志的不可变事件构建的。在这种情况下，事件存储仅支持追加，不鼓励甚至禁止更新或删除操作。事件旨在反映在应用程序级所发生的事情，而不是低级别的状态更改。

事件溯源是一种强大的数据建模技术：从应用程序的角度来看，将用户的行为记录为不可变的事件更有意义，而不是记录这些行为对可变数据库的影响。事件溯源使得随着时间的推移去演化应用程序变得更加容易，通过更容易理解事情发生的原因，以及防范应用程序错误（请参阅本章后面的“不可变事件的优势”）来帮助调试。

举个例子，“学生取消课程注册”事件以一种中立的方式清楚地表达了一个行为，而副作用“从入学表中删除一个条目，并且一条取消的原因被添加到学生反馈表”则嵌入了许多关于稍后使用数据方式的假设。如果引入新的应用程序功能，例如“将座位提供给等待列表中的下一个人”，则采用事件溯源的思路可以轻松地将其集成到现有事件上。

事件溯源类似于编年史数据模型<sup>[45]</sup>，事件日志和星型模式中的事实表也有类似之处（请参阅第3章的“星形与雪花形分析模式”）。

目前，业界已经开发了专门的数据库如Event Store<sup>[46]</sup>来支持使用事件溯源的应用程序，但通常这种方法独立于任何特定的工具。传统的数据库或基于日志的消息代理也可以用来构建这种风格的应用程序。





## 从事件日志导出当前状态

事件日志本身并不是很有用，因为用户通常期望看到系统的当前状态，而不是修改的历史记录。例如，在购物网站上，用户期望能够看到他们购物车的当前内容，而不是对购物车里曾经的修改历史（追加式）。

因此，使用事件溯源的应用程序需要记录事件的日志（表示写入系统的数据），并将其转换为适合向用户显示的状态（从系统读取数据的方式<sup>[47]</sup>）。这种转换可以使用任意的逻辑，但它应该是确定性的，以便可以再次运行它并从事件日志中派生相同的应用程序状态。

与变更数据捕获一样，replay事件日志能够重建系统的当前状态。但是，日志压缩需要以不同的方式处理：

- 用于更新记录的CDC事件通常包含记录的全部新版本，因此key的当前值完全由该key的最近事件确定，并且日志压缩可以丢弃相同key之前的事件。
- 另一方面，使用事件溯源在更高的层次上对事件建模：事件通常用来表达用户行为的意图，而不是一种对行为结果进行相应状态更新的机制。在这种情况下，后来的事件通常不会覆盖以前的事件，所以需要事件的完整历史来重建最终状态。日志压缩不可能以相同的方式进行。

使用事件溯源的应用程序通常有一些机制来保存从导出的当前状态的快照，因此它们不需要重复处理全部的日志。但是，这只是一个性能优化，以加快读取速度和崩溃恢复。其目的是系统能够永久存储所有的原始事件，并在需要时重新处理完整的事件日志。我们会在本章后面的“不可变性的限制”中讨论这个假设。

## 命令和事件

事件溯源的哲学是小心的区分事件和命令<sup>[48]</sup>。当来自用户的请求第一次到达时，它最初是一个命令：此时它可能仍然会失败，例如因为违反了某些完整性条件。应用程序必须首先验证它是否可以执行该命令。如果验证成功并且命令被接受，它将变成一个持久且不可变的事件。

例如，如果用户试图注册某个特定的用户名，或在飞机上或剧院中预订座位，则应用程序需要检查用户名或座位尚未占用（之前在第9章的“支持容错的共识”中讨论过这个例子）。当检查成功时，应用程序可以生成一个事件来指出特定的用户名被特定的用户ID注册了，或者特定的座位已预留给特定的某个客户了。





在事件发生的时候，它成为了事实。即使客户之后决定更改或取消预订，他们以前曾为某个特定的座位进行预订依然是不争的事实，而更改或取消是稍后添加的单独事件。

不允许事件流的消费者拒绝事件：当消费者看到事件时，它已经是日志中不可变的部分，并且可能已经被其他消费者看到。因此，任何命令的验证都需要在它成为事件之前同步发生，例如，通过使用能够原子地验证命令并发布事件的可序列化事务。

或者，预订座位的用户请求可以分成两个事件：第一个是暂时预约，第二个是确认预约后的单独确认事件（如第9章的“采用全序广播实现线性化存储”中所述）。这种划分允许验证过程异步进行。

## 状态，流与不可变性

在第10章介绍了批处理受益于其输入文件的不变性，所以可以在现有的输入文件上运行实验处理作业，而不用担心损坏它们。这种不变性原则也是使得事件溯源和变更数据捕获如此强大的原因。

我们通常将数据库看成是用来存储应用程序当前状态的，这种表示法针对读取进行了优化，并且通常对于查询服务来说是最方便的。状态的本质是它会发生变化，所以数据库支持更新、删除以及插入数据。这如何符合不变性？

每当状态改变，该状态就反映了随着时间推移而变化的事件的结果。例如，当前可用座位列表是已经处理的预订的结果，当前账户余额是账户中的信用和借记的结果，Web服务器的响应时间图是所有发生的Web请求的单个响应时间的集合。

无论状态如何变化，总会有一系列事件导致这些变化。不管事件已经结束或者尚在进行中，事件发生是不争的事实。关键的思路是可变状态和不变事件的追加日志不相互矛盾：它们是同一枚硬币的两面。所有变化的日志，更新日志，代表了随着时间的推移状态的演变。

如果你擅长数学，你可能会说应用状态是事件流对时间的积分得到的，而变化流是状态对时间的求导得到的，如图11-6所示<sup>[49-51]</sup>。虽然这个比喻有一定的局限性（例如，状态的二阶导数似乎没有意义），但可以帮助进一步认知数据。







$$state(now) = \int_{t=0}^{now} stream(t) dt \qquad stream(t) = \frac{d state(t)}{dt}$$

图11-6：当前应用程序状态和事件流之间的关系

如果持久化保存更新日志，不过是实现状态可重现效果。如果认为事件的日志是记录系统，并且从它派生出任何可变状态，那么就更容易推断通过系统的数据流。正如Pat Helland所说<sup>[52]</sup>：

事务日志记录了对数据库所做的所有更改。高速追加是更改日志的唯一方法。从这个角度来看，数据库的内容保存了日志中最新记录值的缓存。日志是事实。数据库是日志子集的缓存。该缓存子集恰好是来自日志的每个记录和索引值的最新值。

日志压缩（如本章前面的“日志压缩”中所述）则是链接日志与数据库区别的一种方式。它仅保留每条记录的最新版本，并丢弃被覆盖的版本。

## 不变事件的优势

数据库中的不变性是一个古老的想法。例如，几个世纪以来，会计师一直在财务簿中使用不变性。交易发生时，它被记录在仅追加的账本中，这本质上是描述货币、商品或服务转手的事件日志。诸如损益或资产负债表之类的账户，则是从账本交易中累加而来<sup>[53]</sup>。

如果发生错误，会计师不会删除或更改账本中的错误交易，而是增加另一笔交易来弥补错误，例如退还不正确的费用。不正确的交易将永远保留在分类账中，出于审计原因这可能很重要。如果从不正确的账本中得出的错误数字已经公布，那么在下一个会计周期会有对应的更正数字。这个过程在会计领域可以说是司空见惯<sup>[54]</sup>。

尽管这种可审计性在金融系统中尤其重要，但对于其他许多不受这种严格规定的系统也是有好处的。如第10章“批处理输出的哲学”中所述，如果意外地部署了将错误数据写入数据库的错误代码，并且代码能够破坏性地覆盖数据，恢复将变得更加的困难。通过不可变事件的追加日志，诊断问题和从问题中恢复就要容易得多。

不可变的事件还会捕获更多的信息，而不仅仅是当前的状态。比如在购物网站，顾客可以将商品添加到购物车，然后再将其移除。从订单执行的角度，虽然第二个事件抵消了第一个事件，但是出于分析的目的，知道客户考虑过某个特定的商品，之后却决



定不购买也很有用。也许他们会选择在未来购买，或者找到了替代品。事件日志会记录所有这些信息，而对于数据库，当他们从购物车中删除时，数据库也删除了相关记录<sup>[42]</sup>。

## 相同的事件日志派生多个视图

此外，通过从不变事件日志中分离可变状态，可以从相同的事件日志派生出多个面向读取的表示方式。这就像有多个消费者的流一样（见图11-5）。例如，分析数据库Druid使用这种方法直接从Kafka获取事件<sup>[55]</sup>，Pistachio是一个分布式的键值存储，使用Kafka作为提交日志<sup>[57]</sup>，Kafka Connect sinks可以将来自Kafka的数据导出到各种不同的数据库和索引中<sup>[41]</sup>。对于许多其他存储和索引系统（如搜索服务器）来说，类似地从分布式日志中获取输入也是有意义的（请参阅本章前面的“保持系统同步”）。

从事件日志到数据库有一个明确的转换步骤，可以更加容易地随时间来演进应用程序：如果想要引入一个新的方式呈现现有数据，可以使用事件日志来构建一个单独针对新功能的读取优化视图，并与现有的系统一起运行，而不需要修改它们。同时运行旧系统和新系统通常比在现有系统中执行复杂的模式迁移更容易。一旦旧的系统不再需要，就可以简单地关闭它并回收它的资源<sup>[47,57]</sup>。

如果不必担心如何去查询和访问数据，那么存储数据通常是非常简单的。模式设计、索引和存储引擎的许多复杂性多是源于希望支持某些查询和访问模式（参看第3章）。因此，将数据写入形式与读取形式分开，并允许多个不同的读取视图，可以获得很大的灵活性。这个想法有时被称为命令查询责任分离（Command Query Responsibility Segregation, CQRS）<sup>[42,58,59]</sup>。

数据库和模式设计的传统方法是基于数据查询必须与数据写入的形式相同这一谬误。如果可以将数据从写优化的事件日志转换为读优化的应用程序状态，有关规范化和非规范化的争论（请参阅第2章的“多对一与多对多关系”）则会变得无关紧要：由于转换过程提供了响应机制使其与源事件日志保持一致，因此在读优化的视图中对数据进行反规范化处理是完全合理的。

在第1章的“描述负载”中，我们讨论了Twitter的主页时间线，一个特定用户正在关注的人最近写的推文的缓存（很像邮箱）。它是读优化的另一个例子：因为你的推文在所有关注你的人的时间线上都是重复的，所以主页时间线是高度非规范化的。然而，扇出服务可以保持重复状态与新的推文和新的关注之间的同步，从而确保对状态重复的可管理性。



## 并发控制

事件捕获和变更数据捕获的最大缺点是事件日志的消费者通常是异步的，所以用户可能会写入日志，然后从日志派生的视图读取，却发现这些写操作还没有反映在读取视图中。我们在第5章的“读自己的写”中讨论了这个问题和可能的解决方案。

一种解决方案是同步执行读取视图的更新，并将事件追加到日志中。这需要一个事务来将写入操作合并到一个原子单元中，所以要么需要将事件日志和读取视图保存在同一个存储系统中，要么需要跨不同系统的分布式事务。或者，可以使用在第9章“采用全序广播实现线性化存储”中讨论的方法。

另一方面，从事件日志导出当前状态也简化了并发控制。对于多对象事务的大部分需求（请参阅第7章的“单对象与多对象事务操作”）源自单个用户需要在不同地方改变数据的操作。通过事件溯源，可以设计一个事件，使其成为用户操作的独立描述。用户操作只需要在一个地方进行一次写操作，即将事件追加到日志中，这很容易使其原子化。

如果以相同的方式对事件日志和应用程序状态进行分区（例如，处理分区3中客户的事件仅需要更新应用程序状态的分区3），则简单的单线程日志消费者不需要对写操作进行并发控制——通过构造，它一次只处理一个事件（另请参阅第7章“实际的串行执行”）。该日志通过在分区中定义事件的串行顺序来消除并发的不确定性<sup>[24]</sup>。如果一个事件涉及多个状态分区，那么需要做更多的工作，将在第12章讨论。

## 不变性的限制

许多不使用事件溯源模型的系统也依赖于不变性：各种数据库在内部使用不可变的数据结构或者多版本数据来支持时间点快照（请参阅第7章“索引与快照隔离”）。诸如 Git、Mercurial 和 Fossil 等版本控制系统也依赖于不可变的数据来保存文件的版本历史记录。

在多大程度上，永远保存所有变化的历史记录是可行的？答案取决于数据集的变化情况。一些工作负载主要是添加数据，很少更新或删除数据，因此很容易支持不变性。其他工作负载在较小的数据集上有较高的更新和删除率，此时，不变的历史数据可能变得过于庞大，碎片化也可能成为一个问题，并且压缩和垃圾回收的性能对于运维的健壮性变得至关重要<sup>[60,61]</sup>。

除了性能之外，还可能在某种情况下，由于管理方面的原因需要删除数据，尽管这些数据都是不可变的。例如，隐私条例可能要求在用户关闭账户后删除他们的个人信息，数据保护法规可能要求删除错误的信息，或者可能需要控制敏感信息的意外泄露。



在这种情况下，仅仅在日志中追加另一个事件来指示先前的数据被视为已删除是明显不够的。实际上，你是想重写历史数据并且假装数据从未写入。例如，Datomic称该特性为*excision*（切除）<sup>[62]</sup>，而Fossil版本控制系统也有一个类似的概念，叫作*shunning*（回避）<sup>[63]</sup>。

真正的删除数据反而会非常困难[64]，这是因为数据副本可能在很多地方都有。例如，存储引擎、文件系统和固态硬盘通常会写入一个新的地址而不是覆盖原地址<sup>[52]</sup>，而备份通常是不可改变的，以防止意外被删除或破坏。删除更多的是“使检索数据更加困难”之意，而不是“使检索数据彻底不可能”。尽管如此，有时必须尝试一下，如将在第12章的“立法与自律”中看到的那样。

## 流处理

本章到目前为止，已经讨论了流的来源（用户活动事件，传感器以及数据库的写操作），讨论了流是如何传输的（通过直接消息传递，通过消息代理和事件日志）。

接下来需要讨论的是，有了流之后，可以用它来做什么，即怎么处理它。一般来说，有三种选择：

1. 可以将事件中的数据写入数据库、缓存、搜索索引或者类似的存储系统，然后被其他客户端查询。如图11-5所示，这是保持数据库与系统其他部分同步的一种好方法，特别是在消费者是写入数据库的唯一客户端的情况下。写入存储系统的流等效于在第10章“批处理工作流的输出”讨论的内容。
2. 可以通过某种方式将事件推送给用户，例如通过发送电子邮件警报或推送通知，或者将事件以流的方式传输到实时仪表板进行可视化。在这种情况下，人是流的最终消费者。
3. 可以处理一个或多个输入流以产生一个或多个输出流。数据流可能会先经过由几个这样的处理阶段组成的流水线，最终在输出端结束（选项1或选项2）。

在本章的余下部分，我们将讨论选项3：处理流以产生其他派生流。处理流的代码逻辑被称为操作或者作业。它与第10章中讨论过的UNIX进程和MapReduce作业密切相关，并且数据流的模式是类似的：流处理器以只读的方式接收输入流，并以仅追加方式将处理输出写入新的位置。

流处理器中的分区和并行化模式也非常类似于第10章中介绍的MapReduce和数据流引擎，因此不在这里重复介绍。基本的映射操作（如转换和过滤记录）也是一样的。

流与批量作业的一个关键区别是，流不会结束。这种差别有很多含义：正如本章开始部分所讨论的，排序对无界数据集没有意义，因此不能使用排序合并join（请参阅第10章“reduce端的join与分组”）。容错机制也必须改变：对于已经运行了几分钟的批处理作业，可以简单地从头开始重新启动失败的任务，但是对于已经运行好几年的流处理作业，在崩溃之后重新开始几乎不可行。

## 流处理的适用场景

流处理长期以来一直被用于监控目的，即希望在发生某些特定事件时收到警报。例如：

- 欺诈检测系统需要确定信用卡的使用方式是否发生了意外的变化，如果信用卡看起来可能已经被盗，就要将其冻结。
- 交易系统需要检查金融市场的价格变化，并根据指定的规则进行交易。
- 制造系统需要监控工厂中机器的状态，在出现故障时快速识别问题。
- 军事和情报系统需要追踪潜在的侵略者的活动，并在有迹象表明发生袭击时发出警报。

这些类型的应用需要相当复杂的模式匹配和相关性。然而，流处理的其他用途也随着时间的推移而出现了。本节，我们将简要比较一下这些应用。

## 复杂事件处理

复杂事件处理（Complex Event Processing，CEP）是20世纪90年代为分析事件流而发展的一种方法，尤其适用需要搜索特定的事件模式<sup>[65,66]</sup>。与正则表达式支持在字符串中搜索特定字符模式的方式类似，CEP允许指定规则，从而可以在流中搜索特定模式的事件。

CEP系统通常使用像SQL这样的高级声明式查询语言或图形用户界面，来描述应该检测到的事件模式。这些查询提交给一个处理引擎，该引擎使用输入流并在内部维护匹配所需的状态机。当发现匹配时，引擎产生一个复杂的事件，这个事件包括检测到的事件模式的细节信息<sup>[67]</sup>。

在这些系统中，查询和数据之间的关系与普通数据库相比正好相反。通常情况下，数据库会持久存储数据并将查询视为暂时的：当查询到来时，数据库搜索与查询匹配的数据，然后在查询完成时忘记它。CEP引擎则反转了这些角色：查询是长期存储的，来自输入流的事件不断流过他们以匹配事件模式<sup>[68]</sup>。

CEP的实现包括Esper<sup>[69]</sup>、IBM Info Sphere Streams<sup>[70]</sup>、Apama、TIBCO StreamBase和SQLstream。像Samza这样的分布式流处理器也对流支持声明式SQL查询<sup>[71]</sup>。

## 流分析

使用流处理的另一个领域是对流进行分析。CEP和流分析之间的界限有些模糊，但作为一般规则，分析往往不太关心找到特定的事件序列，而更多地面向大量事件的累计效果和统计指标，例如：

- 测量某种类型事件的速率（每个时间间隔发生的频率）。
- 计算一段时间内某个值的滚动平均值。
- 将当前的统计数据与以前的时间间隔进行比较（例如，检测趋势或提醒与上周同一时间相比异常高或低的指标）。

这些统计信息通常是在固定的时间间隔内进行计算的，例如，过去5min内每秒对服务的平均查询次数是多少，以及在此期间的第99个百分位的响应时间；几分钟内的平均值可以消除从一秒内无关紧要的波动，同时还能及时了解流量的最新变化。聚合操作的时间间隔称为窗口，我们将在本章后面“关于时间的推理”中更详细地讨论窗口。

流分析系统有时使用概率算法，比如用于设置成员关系的布隆过滤器（在第3章的“性能优化”中介绍过），用于基数估计的HyperLogLog，以及各种百分比估值计算方法（请参阅第1章“实践中的百分位数”）。概率算法只能产生近似的结果，但其优点是在流处理器中所需的内存明显少于精确算法。近似算法的使用有时会导致人们认为流处理系统总是失真和不够精确，但这种观点是错误的：流处理本身并没有任何固有的近似处理，概率算法仅仅是一种优化<sup>[73]</sup>。

许多开源的分布式流处理框架在设计时都考虑了对分析的支持，例如，Apache Storm、Spark Streaming、Flink、Concord、Samza和Kafka Streams<sup>[74]</sup>。类似的，还有一些托管服务如Google Cloud Dataflow和Azure Stream Analytics。

## 维护物化视图

在本章前面“数据库和数据流”看到，可以使用数据库更改流来保持派生数据系统（如缓存、搜索索引和数据仓库等）与源数据库之间的同步。可以将这些示例视为一种维护物化视图的例子（请参阅第3章“聚合：数据立方体与物化视图”）：对某个数据集导出一个特定的视图以便高效查询，并在底层数据更改时自动更新该导出视图<sup>[50]</sup>。



同样，在事件溯源中，应用状态通过应用事件日志来维护，这里的应用状态也是一种物化视图。与流分析场景不同，仅考虑某个时间窗口内的事件通常还不够：除了那些可能被日志压缩丢弃的事件（请参阅本章前面“日志压缩”），构建物化视图可能需要任意时间段内的所有事件。实际上，需要的是一个可以延伸到开始时间的足够长的窗口。

原则上，任何流处理器都可以用于物化视图维护，尽管永久维护事件的需求与面向分析的框架的假设背道而驰，这些框架主要在有限持续时间的窗口上运行。Samza和Kafka Streams支持这种用法，基于Kafka的日志压缩功能<sup>[75]</sup>。

## 在流上搜索

CEP通常搜索包含多个事件的特定模式，除了CEP，有时还需要基于一些复杂条件（例如全文搜索查询）来搜索单个事件。

举个例子，媒体监控服务订阅来自媒体机构的新闻文章或者公告，并支持搜索关于公司、产品或感兴趣主题的相关新闻。这是通过预先制定一个搜索查询来完成的，然后不断地将新闻流与这个查询进行匹配。在一些网站上也有类似的功能：例如，房地产网站的用户需要当市场上出现符合其搜索条件的新房产时被及时通知。Elasticsearch<sup>[76]</sup>的过滤器功能是实现这种流式搜索的一种方式。

传统的搜索引擎首先索引文档，然后在索引上运行查询。相比之下，搜索流则是反过来：查询条件先保存下来，所有文档流过查询条件，就像CEP一样。最极端的情况，可以针对每个查询来测试每个文档，但是如果有大量的查询，这可能会非常缓慢。为了优化，可以对查询和文档都进行索引，从而缩小匹配的查询集合<sup>[77]</sup>。

## 消息传递和RPC

在第4章“基于消息传递的数据流”中，我们讨论了消息传递系统作为RPC的替代方法，即作为通信服务的一种机制，例如actor模型中所使用的机制。虽然这些系统也是基于消息和事件，但是我们通常并不把他们归类为流处理系统：

- Actor框架主要是管理通信模块的并发和分布式执行的机制，而流处理主要是数据管理技术。
- Actor之间的交流往往是短暂的，并且是一对一的，而事件日志是持久的、多用户的。
- Actor可以以任意方式进行通信（包括循环请求/响应模式），但流处理器通常设

置在非循环流水线中，其中每个流是一个特定作业的输出，并且从一组定义明确的输入流派生而来。

也就是说，RPC类系统和流处理之间有一些交叉的地方。例如，Apache Storm有一个称为分布式RPC的功能，它支持将用户查询分布到一组同时运行事件处理的节点上，然后这些查询与来自输入流的事件交织在一起，查询结果可以从多节点聚合而来然后返回给用户<sup>[78]</sup>（另请参阅第12章的“多分区数据处理”）。

你也可以使用Actor框架来处理流。然而，许多这样的框架在崩溃的情况下不能保证消息的传递，所以处理不是容错的，除非实现额外的重试逻辑。

## 流的时间问题

流处理系统经常需要和时间打交道，尤其是在用于分析目的时，这些分析通常使用时间窗口，例如“最近五分钟内的平均值”。似乎“最后五分钟”的含义应该明确无误，但不幸的是，这种定义处理起来非常棘手。

在批处理过程中，处理任务会快速处理大量的历史事件。如果需要按时间进行某种分解，批处理需要查看每个事件所嵌入的时间戳。因为处理运行的时间与事件实际发生的时间无关，所以查看运行批处理的机器的系统时钟没有意义。

批处理可以在几分钟内读取一年的历史事件；在大多数情况下，关注的是一年的历史事件，而不是几分钟的处理过程。而且，在事件中使用时间戳可以使得处理过程是确定性的：基于同一个输入，再次运行相同的处理过程可以得到相同的结果（请参阅第9章的“故障容错”）。

另一方面，许多流处理框架使用处理节点上的本地系统时钟（处理时间）来确定窗口<sup>[79]</sup>。这种方法的优点是简单，当事件发生和事件处理之间的间隔可以忽略不计，这种方式也是合理的。然而，如果存在显著的处理滞后，即处理可能比事件实际发生的时间明显要晚，则该方法不再有效。

## 事件时间与处理时间

发生事件处理滞后的原因有很多，例如排队、网络故障（请参阅第8章“不可靠的网络”）、导致消息代理或处理器中出现竞争的性能问题，重新启动流的消费者，重新处理过去的事件（请参阅本章前面“重新处理消息”，例如从错误中恢复或修复了代码的错误）。



而且，消息延迟还可能导致消息的不可预知的排序。例如，假设用户首先发出一个Web请求（由Web服务器A处理），然后发出第二个请求（由服务器B处理）。A和B发出描述它们所处理的请求的事件，但是B的事件在A的事件发生之前到达消息代理。现在，流处理器将首先看到B事件，然后看到A事件，即使它们实际上是以相反的顺序发生的。

如果找一个类比的话，可以考虑一下“星球大战”这部电影<sup>注2</sup>：第四部于1977年上映，1980年的第五部，1983年的第六部，之后分别于1999年，2002年和2005年上映第一部，第二部，第三部，以及2015年的第七部<sup>[80]</sup>。如果以上映顺序来观看电影，则观看（“处理”）电影的顺序与它们叙述的事件顺序就是不一致的（每部的编号就像事件时间戳一样，而观看电影的日期就是处理时间）。人类能够应付这样的不连续性，但是流处理算法则需要专门的代码处理，以适应这样的时间和排序问题。

混淆事件时间与处理时间会导致错误的结果。例如，假设有一个流处理操作来测量请求频率（计算每秒的请求数）。如果重新部署了该流处理系统，则中间可能会关闭一分钟，在重启后继续处理积压的事件。如果根据处理时间来衡量请求频率，那么看起来好像在处理积压事件时突然出现了异常的请求高峰，而事实上则是请求频率一直是稳定的（见图11-7）。

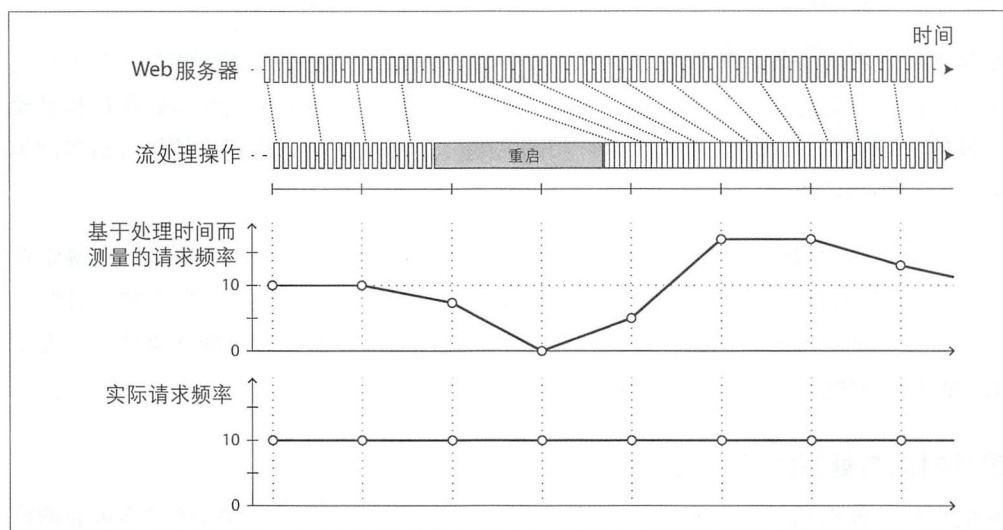


图11-7：由于处理速率的波动，基于处理时间的窗口会引入结果异常

注2：感谢Flink社区的Kostas Kloudas提出星球大战电影的这个类比例子。







## 了解什么时候准备就绪

如果基于事件发生时间而定义窗口，面临一个棘手的问题是，你无法确定什么时候能收到特定窗口内的所有事件，或者是否还有一些事件尚未到来。

例如，假设将事件分组成一分钟的窗口，以便可以统计每分钟的请求数。你已经计算了一些事件，这些事件的时间戳在一小时的第37分钟，时间继续向前移动，现在大部分事件都到了第38分钟和第39分钟。什么时候该宣布已经完成了第37分钟窗口的处理，可以输出计数器值？

在一段时间没有看到任何新的事件之后，可以认为超时并宣布关闭该窗口，但由于网络中断而延迟，仍然可能发生某些事件其实还缓存在另一台计算机上。需要能够处理在窗口已经声明完成后才到达的这样的滞后事件。大体上，有两个选择<sup>[1]</sup>：

1. 忽略这些滞后的事件，因为他们在正常情况下可能只是一小部分事件。可以将丢弃事件的数量作为度量标准进行跟踪，当出现丢弃大量数据时发出警报。
2. 发布一个更正：针对滞后事件的一个更新值。可能还需要收回以前的输出。

在某些情况下，可以使用一个特殊的消息来表明，“从现在开始，不会有比t更早的时间戳的消息”，消费者可以使用它来触发窗口处理<sup>[81]</sup>。但是，如果不同机器上的多个生产者正在生成事件，每个事件都有自己的最小时间戳阈值，则消费者需要分别跟踪每个生产者。在这种情况下添加和删除生产者变得比较麻烦。

## 你用谁的时钟？

当事件可能在系统中的多个点缓冲时，为事件分配时间戳就比较困难。例如，考虑向服务器报告使用率事件的移动应用。该应用可能会在设备处于离线状态时使用，在这种情况下，它会在设备本地缓存事件，并在下一次有可用的网络连接（可能是几小时甚至几天）时将其发送到服务器上。对于这个流的任何一个消费者来说，这些事件都是极其滞后的。

这时，根据移动设备的本地时钟，事件的时间戳实际上指的是发生交互时的时间。然而，用户控制的设备上的时钟通常是不可信的，它可能会被意外或故意设置为错误的时间（请参阅第7章“时钟同步与精度”）。服务器收到事件的时间（根据服务器的时钟）更可能是准确的，因为服务器在你的控制之下，但是在描述用户交互方面意义就不大了。

为了调整不正确的设备时钟，一种方法是记录三个时间戳<sup>[82]</sup>：





- 根据设备的时钟，记录事件发生的时间。
- 根据设备的时钟，记录将事件发送到服务器的时间。
- 根据服务器时钟，记录服务器收到事件的时间。

通过从第三个时间戳中减去第二个时间戳，可以估计出设备时钟和服务器时钟之间的偏移量（假设网络延迟与所需的时间戳精度相比可以忽略不计）。然后，可以将该偏移量应用于事件时间戳，从而估计事件实际发生的真实时间（假设设备时钟偏移量在事件发生的时间与发送到服务器的时间之间没有发生变化）。

这个问题并不是流处理独有的，批处理也会在时间推理方面遇到相同的问题。只不过在流处理环境中更加明显，因为我们更能注意到时间的变化。

## 窗口类型

一旦明确了如何确定事件的时间戳，下一步就是决定如何定义时间段即窗口了。然后窗口就可用于聚合分析，例如计数事件，或计算窗口内的所有值的平均值。有以下几种常见的窗口类型<sup>[79,83]</sup>：

### 轮转窗口

翻滚窗口的长度是固定的，每个事件都属于一个窗口。例如，如果有一个1分钟的翻滚窗口，则所有时间戳在10:03:00和10:03:59之间的事件都会被分组到一个窗口中，10:04:00和10:04:59之间的事件被分组到下一个窗口，依此类推。可以通过获取每个事件时间戳并将其四舍五入到最近的分钟来确定它所属的窗口，从而实现1分钟的轮转窗口。

### 跳跃窗口

跳跃窗口也具有固定长度，但允许窗口重叠以提供一些平滑过渡。例如，一个5分钟窗口，设定跳跃（hop）值为1分钟，则包含10:03:00至10:07:59之间的事件，而下一个窗口将包括10:04:00至10:08:59之间的事件，依此类推。可以通过首先计算1分钟滚动窗口，然后聚合几个相邻的窗口来实现跳跃窗口。

### 滑动窗口

滑动窗口包含在彼此的某个间隔内发生的所有事件。例如，一个5分钟的滑动窗口将包括10:03:39和10:08:12的事件，因为它们相距不到5分钟（注意，轮转和跳跃的5分钟窗口不会将这两个事件置于相同的窗口，因为它们都基于固定的边界）。滑动窗口可以通过保留按时间排序的事件缓冲区并且在从窗口过期时移除旧事件来实现。





## 会话窗口

与其他窗口类型不同，会话窗口没有固定的持续时间。相反，它是通过将同一用户在时间上紧密相关的所有事件分组在一起而定义的，一旦用户在一段时间内处于非活动状态（例如，如果30分钟内没有事件），则窗口结束。会话分析是网站分析中常见的一种需求（参阅第10章“分组”）。

## 流式join

在第10章中，我们讨论了批处理作业如何通过主键来join数据集，以及这种join如何构成数据管道。由于流处理将数据管道推广为对无界数据集的增量处理，因此对流进行join的需求也完全适用。

但是，新事件随时可能出现在流中，这个事实使得流式join比批处理作业更具挑战性。为了更好地理解，我们来区分三种不同类型的join：流和流join，流和表join和表和表join<sup>[84]</sup>。在下面的章节中，我们将通过例子来逐一说明。

### 流和流join（窗口join）

假设某网站支持搜索功能，并且想要检测网址搜索的最新趋势。每次有人输入搜索查询时，都会记录包含查询和返回结果的事件。每当有人单击其中一个搜索结果时，就会记录另一个单击的事件。为了计算搜索结果中每个网址的单击率，需要将搜索操作和单击操作的事件组合在一起，这些事件通过具有相同的会话ID进行join。广告系统也需要类似的分析<sup>[85]</sup>。

如果用户并不相信搜索结果，可能永远不会发生单击事件，即使发生了这样的情况，搜索和单击之间的时间也可能存在很大间隔：在许多情况下可能是几秒钟，但是可能长达几天或几周（如果用户运行完搜索后忘记了浏览器的这个选项卡，过了些时间后又打开这个选项卡，再单击这个搜索结果）。由于网络延迟也是变化的，单击事件甚至可能在搜索事件之前到达。这时可以定义合适的join窗口，例如，如果搜索和单击的间隔不超过一小时，则可以join它们。

请注意，在单击事件中嵌入搜索的细节并不等同于join事件：这样做只会告诉你用户单击了什么搜索结果，而无法告诉你用户没单击时的情况。评估搜索质量需要准确的单击率，因此需要搜索事件和单击事件两个事件源。

为了实现这种类型的join，流处理操作需要维护状态：例如，在最后一小时发生的所有事件（基于会话ID索引）。无论什么时候发生搜索事件或单击事件，都会将其添加







到适当的索引，并且流处理还检查另一个索引，以查看同一个会话ID的另一个事件是否已到达。如果有匹配的事件，则发出一个派生事件来表明哪个搜索结果发生了单击。如果搜索事件过期而没有匹配到单击事件，则发出一个派生事件表明哪些搜索结果未被单击。

## 流和表join

在第10章的“示例：用户行为事件的分析”（见图10-2）中，我们介绍了join两个数据集的批处理任务例子：一组用户活动事件和一个用户资料。将用户活动事件视为流，并在流处理器中连续执行相同的join是比较自然的处理方法：输入是包含用户ID的活动事件流，输出是在用户ID上追加了用户的资料信息的活动事件流。这个过程有时被称为使用来自数据库的信息丰富活动事件。

要执行此join，流处理过程需要一次查看一个活动事件，在数据库中查找事件的用户ID，然后将该概要信息添加到活动事件中。数据库查询可以通过查询远程数据库来实现。但是，正如在第10章“示例：分析用户活动事件”中讨论的，此类远程查询可能很慢并且有可能导致数据库过载<sup>[75]</sup>。

另一种方法是将数据库副本加载到流处理器中，以便在本地进行查询而无需经过网络往返。这种技术与第10章“Map端的join操作”中所讨论的哈希join非常相似：本地副本可能是位于内存中的哈希表（如果数据比较小），或者是本地磁盘上的索引。

与批处理任务的区别在于，批处理任务使用数据库的时间点快照作为输入，而流处理器是长时间运行的，并且数据库的内容可能随时间而改变，所以流处理器数据库的本地副本需要保持最新。这个问题可以通过变更数据捕获手段来解决：流处理器可以订阅用户资料数据库的更新日志以及活动事件流。在创建或修改资料时，流处理器会更新其本地副本。因此，我们获得两个流之间的join。

流和表join实际上非常类似流和流join。它们最大的区别在于，对于表的更新日志流，join使用一个可以回溯到“开始时间”（一个在概念上是无限的窗口）的窗口，而新版本的记录会覆盖旧的记录。对于流输入，join可能根本无法维护这样一个窗口。

## 表和表join（物化视图维护）

考虑在第1章“描述负载”时Twitter时间线的例子。当用户想要查看其主页时间线时，循环遍历用户关注的所有人、查找他们最近的推文并将其合并，这个操作代价很高。

相反，我们需要一个时间线缓存：这是一种基于每个用户的“收件箱”，在发送推文





的时候就将其写入其中，因此读取时间线就只需要一次查询。实现和维护这个缓存需要以下事件处理：

- 当用户 $u$ 发送新的推文时，它被添加到每个关注你的用户的时间线上。
- 当用户删除推文时，会从所有用户的时间线中删除。
- 当用户 $u_1$ 开始关注用户 $u_2$ ， $u_2$ 最近的推文被添加到 $u_1$ 的时间线上。
- 当用户 $u_1$ 取消关注用户 $u_2$ 时，由 $u_2$ 发出的推文将从 $u_1$ 的时间线中移除。

要在流处理中实现这种缓存维护，需要用于推文（发送和删除）和关注关系（关注和取消关注）的事件流。流过程需要维护一个包含每个用户的关注者集合的数据库，以便知道当一个新的推文到达时需要更新哪些时间线<sup>[86]</sup>。

理解这个流处理过程的另一种方法是，它维护一个两个表（推文和关注者）join查询的物化视图，如下所示：

```
SELECT follows.follower_id AS timeline_id,
       array_agg(tweets.* ORDER BY tweets.timestamp DESC)
FROM tweets
JOIN follows ON follows.followee_id = tweets.sender_id
GROUP BY follows.follower_id
```

流的join直接对应于该查询中的表的join。时间线实际上是这个查询结果的缓存，每当底层表发生变化时就进行更新<sup>注3</sup>。

## join的时间依赖性

这里描述的三种类型的join（流-流，流-表和表-表）有很多相同之处：它们都需要流处理系统根据一个join输入来维护某些状态（例如搜索和单击事件，用户资料或关注列表），并在来自另一个join输入的消息上查询该状态。

维持这个状态的事件的顺序非常重要（它关系到是先关注然后解除关注，或是相反）。在分区日志中，单个分区内的事件排序被保留，但通常在不同的流或分区之间没有排序的保证。

这就产生了一个问题：如果不同流中的事件发生在相近的时间里，它们按照何种顺序

---

注3：如果把流视为表的导数，如图11-6所示，并且把联结看作是两个表 $u \cdot v$ 的乘积，那么会发生一些有趣的事情：物化联结的变化流遵循积分的求导规则  $(u \cdot v)' = u'v + uv'$ 。换句话说，任何推文的变化都与当前的关注者联结在一起，任何关注者的变化都与当前的推文<sup>[49,50]</sup>相联结。





进行处理？在流和表join的示例中，如果用户更新了他们的资料，哪些活动事件会与旧资料（在资料更新之前处理）join，哪些又与新资料join（在资料更新之后处理）？换句话说：如果状态随时间而改变，而join操作需要输入状态信息，那么应该使用什么时间点的状态来join呢<sup>[45]</sup>？

这种时间依赖性可能会在许多地方发生。比如，销售某些商品，需要对发票计算适当的税率，而税率取决于不同的国家或者州、产品类型和销售日期（因为税率会随时变化）。将销售情况与税率表join时，就希望join销售时所对应的税率，而如果正在重新处理历史数据，这可能就与当时的税率不同。

如果跨流的事件排序是不确定的，那么join也变成非确定性的<sup>[87]</sup>，这意味着你不能在相同的输入上重新运行同一作业来得到相同的结果：因为再次运行任务时，输入流上的事件可能以不同的方式交叉在一起。

在数据仓库中，这个问题被称为缓慢变化的维度（Slowly Changing Dimension，SCD），通常通过对特定版本的join记录赋予唯一的标识符来解决：例如，每当税率改变时，就赋予一个新的标识符，而发票也包括销售时的税率标识符<sup>[88,89]</sup>。这种方式使join操作具有确定性，但是由于表中所有不同版本的记录都需要保留，最后日志压缩几乎无法有效工作。

## 流处理的容错

在本章的最后部分，我们来考虑一下流处理器如何容错。在第10章中看到，批处理框架可以比较容易实现容错：如果MapReduce作业中的任务失败，可以简单地在另一台机器上重新启动，并丢弃失败任务的输出。这种透明的重试是可能的，因为输入文件是不可变的，每个任务都将其输出写到HDFS上的单独文件，并且输出仅在任务成功完成时可见。

特别是，批处理容错方法可以确保批处理作业的输出与没有出错时的最终结果相同，即使事实上某些任务失败了。看起来好像每个输入记录都恰好被处理了一次，没有记录被遗漏，也没有记录被处理两次。尽管重新启动任务意味着实际上可能会多次处理记录，但在输出中的效果看起来和只处理过一次一样。这个原则被称为恰好一次语义，虽然有效一次应该更为准确<sup>[90]</sup>。

在流处理过程中也会出现同样的容错问题，但是处理起来并不那么简单：在使输出结果可见之前等待某个任务完成是不可行的，由于流是无限的，因此几乎永远无法完成这个任务。







## 微批处理和校验点

一种解决方案是将流分解成多个小块，并像小型批处理一样处理每个块。这种方法被称为微批处理，它已经用于Spark Streaming<sup>[91]</sup>。批处理大小通常约为1s，这是一个性能折中的考虑：较小的批处理会导致更大的调度和协调开销，而较大的批处理意味着流处理器的结果需要更长的延迟才能可见。

微批处理隐含地设置了与批处理大小相等的轮转窗口（基于处理时间而不是事件时间）。任何需要更大窗口的作业都需要显式地将状态从一个微批处理保留至下一个微批处理。

Apache Flink中使用了该方法的一个变体，它定期生成状态滚动检查点并将其写入持久化存储<sup>[92,93]</sup>。如果流操作发生崩溃，它可以从最近的检查点重新启动，并丢弃在上一个检查点和崩溃之间生成的所有输出。检查点是由消息流中的barrier触发，类似于微批处理之间的边界，但并不强制特定的窗口大小。

在流处理框架的范围内，微批处理和检查点的方法提供了与批处理一样的恰好一次语义。但是，一旦输出脱离了流处理系统（例如，通过写入数据库，向外部消息代理发送消息或发送电子邮件），框架将无法丢弃失败批处理的输出了。在这种情况下，重启失败的任务将导致发生两次外部可见的副作用，单独使用微批处理或检查点都不足以防范此类问题。

## 重新审视原子提交

在出现故障时，为了看起来实现恰好处理了一次，我们需要确保当且仅当处理成功时，所有输出和副作用才会生效。这包括发送给下游的操作或外部消息传递系统（包括邮件或推送通知）的任何消息，所有数据库的写入，以及对操作状态的更改和任何对输入消息的确认（包括基于日志的消息代理中，向前移动消费者偏移量）。

这些事情要么原子的发生，要么都不发生，但不应该彼此不同步。这种方法听起来很熟悉，因为在第9章“恰好一次消息处理”中分布式事务和两阶段提交的上下文中已经有所讨论。

在第9章中，讨论了分布式事务的传统实现中的问题，比如XA。然而，在更受限制的环境中，有可能有效地实现这样的原子提交方法。Google Cloud Dataflow<sup>[81,92]</sup>和VoltDB<sup>[94]</sup>中使用了这种方法，Apache Kafka<sup>[95,96]</sup>也计划添加类似的功能。与XA不同，这些实现不会尝试跨异构技术提供事务，而是通过在流处理框架中管理状态更改

和消息传递来保持内部事务。事务协议的开销可以通过在单个事务中处理多个输入消息来分摊。

## 幂等性

我们的目标是丢弃任何失败任务的部分输出，以便它们可以安全地重试而不会两次生效。分布式事务是实现这一目标的一种方式，而另一种方式则是依赖幂等性<sup>[97]</sup>。

幂等操作是可以多次执行的操作，并且它与只执行一次操作具有相同的效果。例如，将键-值存储中的键设置为某个固定值是幂等的（再次写入该值会覆盖已有的相同值），而递增计数器则不满足幂等（再次执行递增意味着该值递增两次）。

即使操作不是天生具备幂等性，往往也可以使用一些额外的元数据使其变得幂等。例如，处理来自Kafka的消息时，每条消息都有一个持久的、单调递增的偏移量。将值写入外部数据库时，可以在该值中包含触发最新写入消息所对应的偏移量。因此，根据该偏移量可以判断是否已经应用了该更新，从而避免重复执行相同的更新操作。

Strom的Trident的状态处理也是基于类似的想法<sup>[78]</sup>。依赖于幂等性意味着基于一些假设：重启失败的任务必须以相同的顺序重新处理相同的消息（基于日志的消息代理是这么做的），处理必须是确定性的，并且没有其他节点并发更新相同的值<sup>[98, 99]</sup>。

当从一个处理节点切换至另一个节点时，可能需要必要的fencing措施（请参阅第8章“主节点与锁”），以防止标记为失败但实际上仍处于活动状态的节点的干扰。尽管存在这些禁忌，幂等操作依然是一种有效的方式来实现恰好一次语义，并且开销很小。

## 故障后重建状态

任何需要状态的流处理，比如基于窗口的聚合（诸如计数器，平均值和直方图）以及表和索引的join操作，都必须确保在故障发生后状态可以恢复。

一种选择是将状态保存在远程存储中并采取复制，然而为每个消息去查询远程数据库可能会很慢。另一种方法是将状态在本地保存，并定期进行复制。之后，当流处理器从故障中恢复时，新任务可以读取副本的状态并且在不丢失数据的情况下恢复处理。

例如，Flink定期对操作状态执行快照，并将他们写入HDFS等持久性存储中<sup>[92, 93]</sup>。Samza和Kafka Streams通过将状态更改发送到具有日志压缩功能的专用Kafka主题来保存状态的副本，这类似于变更数据捕获<sup>[84, 100]</sup>。VoltDB则通过并在多个节点上冗余处理每个输入消息来复制状态（请参阅第7章“实际的串行执行”）。

在某些情况下，甚至可能不需要复制状态，而是从输入流开始重建。例如，如果状态由相当短的窗口中的聚合组成，简单地replay与该窗口相对应的输入事件就可能足够快。如果状态是数据库的本地副本（基于变更数据捕获方式来维护），那么也可以从日志压缩的更改流那里重建数据库（请参阅本章前面“日志压缩”）。

但是，所有这些权衡取决于底层基础架构的性能表现：在某些系统中，网络延迟可能低于磁盘访问延迟，网络带宽可能与磁盘带宽相当。不存在适用于所有情况的理想的折中方案，随着存储和网络技术的发展，本地和远程状态的优点也可能会发生变化。

## 小结

在本章中，我们讨论了事件流，他们的用途以及处理方法。在某些方面，流处理与第10章中讨论的批处理非常相似，但是它是对无限（永不停止）的流而不是固定大小的输入进行持续的处理。从这个角度讲，消息代理和事件日志可以视为流处理系统中文件系统。

我们花了一些时间比较两种类型的消息代理：

### AMQP / JMS风格的消息代理

代理将单个消息分配给消费者，消费者在成功处理后确认每条消息。消息被确认后从代理中删除。这种方法适合作为一种异步RPC（另请参阅第4章“基于消息传递的数据流”），例如，在任务队列中，消息处理的确切顺序并不重要，并且不需要在他们处理完后返回并再次读取旧消息。

### 基于日志的消息代理

代理将分区中的所有消息分配给相同的消费者节点，并始终以相同的顺序发送消息。通过分区机制来实现并行，消费者通过检查他们处理的最后一条消息的偏移量来跟踪进度。代理将消息保存在磁盘上，因此如果有必要，可以回退并重新读取旧消息。

基于日志的方法与在数据库（请参阅第5章）和日志结构化存储引擎（请参阅第3章）中的复制日志相似。我们看到，这种方法特别适合流处理系统，即不断汲取输入流并生成派生状态或派生输出流。

就流的来源而言，我们讨论了几种可能性，包括用户活动事件，提供定期读数的传感器，以及可订阅数据源（如金融市场数据）等可以自然地表示为流。我们还看到，将对数据库的写入视为流也是有用的：可以捕获变更日志（即对数据库所做的所有更改



的历史记录），包括隐式地通过变更数据捕获或显式地通过事件溯源两种方法。通过日志压缩可以保留数据库内容的完整副本。

将数据库表示为流为高效集成系统提供了更多的机会。通过使用变更日志并将其应用到派生系统，可以不断更新搜索索引、缓存和分析系统等派生数据系统。甚至可以全新做起，将所有数据从始至终都视为变更日志，从而构建全新的视图。

将状态保持为流并支持重做消息也是在各种流处理框架中实现基于流的join和容错的基础。我们讨论了流处理的几个目标场景，包括基于事件模式的查询（复杂事件处理），计算窗口聚合（流分析）以及保持派生数据系统处于最新状态（物化视图）等。

接下来讨论了在流处理中有关时间方面的一些考虑，包括处理时间和事件时间之间的区别，以及滞后事件（在认为窗口已关闭之后才到达的事件）问题。

我们区分了流处理过程中可能出现的三种类型的join：

### 流和流join

两个输入流都由活动事件组成，采用join操作用来搜索在特定时间窗口内发生的相关事件。例如，它可以匹配相同用户在30min内采取的两个动作。如果想要在一个流中查找相关事件，则两个join输入可以是相同的流。

### 流和表join

一个输入流由活动事件组成，而另一个则是数据库变更日志。更新日志维护了数据库的本地最新副本。对于每个活动事件，join操作用来查询数据库并输出一个包含更多信息的事件。

### 表和表join

两个输入流都是数据库更新日志。在这种情况下，一方的每一个变化都与另一方的最新状态相join。结果是对两个表之间join的物化视图进行持续的更新。

最后，我们讨论了在流处理器中实现容错和恰好一次语义的技术。与批处理一样，我们需要丢弃所有失败任务的部分输出。然而，由于流处理长时间运行并持续产生输出，所以不能把所有的输出都简单地丢弃掉。相反，基于微批处理、检查点、事务或幂等性写入等，可以实现更细粒度的恢复机制。

## 参考文献

[1] Tyler Akidau, Robert Bradshaw, Craig Chambers, et al.: “The Dataflow Model:

A Practical Approach to Balancing Correctness, Latency, and Cost in Massive-Scale, Unbounded, Out-of-Order Data Processing,” *Proceedings of the VLDB Endowment*, volume 8, number 12, pages 1792-1803, August 2015. doi:10.14778/2824032.2824076.

[2] Harold Abelson, Gerald Jay Sussman, and Julie Sussman: *Structure and Interpretation of Computer Programs*, 2nd edition. MIT Press, 1996. ISBN: 978-0-262-51087-5, available online at [mitpress.mit.edu](http://mitpress.mit.edu).

[3] Patrick Th. Eugster, Pascal A. Felber, Rachid Guerraoui, and Anne-Marie Kermarrec: “The Many Faces of Publish/Subscribe,” *ACM Computing Surveys*, volume 35, number 2, pages 114-131, June 2003. doi:10.1145/857076.857078.

[4] Joseph M. Hellerstein and Michael Stonebraker: *Readings in Database Systems*, 4th edition. MIT Press, 2005. ISBN: 978-0-262-69314-1, available online at [redbook.cs.berkeley.edu](http://redbook.cs.berkeley.edu).

[5] Don Carney, Uğur Çetintemel, Mitch Cherniack, et al.: “Monitoring Streams - A New Class of Data Management Applications,” at *28th International Conference on Very Large Data Bases (VLDB)*, August 2002.

[6] Matthew Sackman: “Pushing Back,” *lshift.net*, May 5, 2016.

[7] Vicent Martí: “Brubeck, a statsd-Compatible Metrics Aggregator,” *githubengineering.com*, June 15, 2015.

[8] Seth Lowenberger: “MoldUDP64 Protocol Specification V 1.00,” *nasdaqtrader.com*, July 2009.

[9] Pieter Hintjens: *ZeroMQ – The Guide*. O’Reilly Media, 2013. ISBN: 978-1-449-33404-8.

[10] Ian Malpass: “Measure Anything, Measure Everything,” *codeascraft.com*, February 15, 2011.

[11] Dieter Plaetinck: “25 Graphite, Grafana and statsd Gotchas,” *blog.raintank.io*, March 3, 2016.

[12] Jeff Lindsay: “Web Hooks to Revolutionize the Web,” *progrium.com*, May 3, 2007.

- [13] Jim N. Gray: “Queues Are Databases,” Microsoft Research Technical Report MSR-TR-95-56, December 1995.
- [14] Mark Hapner, Rich Burrridge, Rahul Sharma, et al.: “JSR-343 Java Message Service(JMS) 2.0 Specification,” *jms-spec.java.net*, March 2013.
- [15] Sanjay Aiyagari, Matthew Arrott, Mark Atwell, et al.: “AMQP: Advanced Message Queuing Protocol Specification,” Version 0-9-1, November 2008.
- [16] “Google Cloud Pub/Sub: A Google-Scale Messaging Service,” *cloud.google.com*, 2016.
- [17] “Apache Kafka 0.9 Documentation,” *kafka.apache.org*, November 2015.
- [18] Jay Kreps, Neha Narkhede, and Jun Rao: “Kafka: A Distributed Messaging System for Log Processing,” at *6th International Workshop on Networking Meets Databases (NetDB)*, June 2011.
- [19] “Amazon Kinesis Streams Developer Guide,” *docs.aws.amazon.com*, April 2016.
- [20] Leigh Stewart and Sijie Guo: “Building DistributedLog: Twitter’s High-Performance Replicated Log Service,” *blog.twitter.com*, September 16, 2015.
- [21] “DistributedLog Documentation,” Twitter, Inc., *distributedlog.io*, May 2016.
- [22] Jay Kreps: “Benchmarking Apache Kafka: 2 Million Writes Per Second (On Three Cheap Machines),” *engineering.linkedin.com*, April 27, 2014.
- [23] Kartik Paramasivam: “How We’re Improving and Advancing Kafka at LinkedIn,” *engineering.linkedin.com*, September 2, 2015.
- [24] Jay Kreps: “The Log: What Every Software Engineer Should Know About Real-Time Data’s Unifying Abstraction,” *engineering.linkedin.com*, December 16, 2013.
- [25] Shirshanka Das, Chavdar Botev, Kapil Surlaker, et al.: “All Aboard the Databus!,” at *3rd ACM Symposium on Cloud Computing (SoCC)*, October 2012.
- [26] Yogeshwer Sharma, Philippe Ajoux, Petchean Ang, et al.: “Wormhole: Reliable Pub-Sub to Support Geo-Replicated Internet Services,” at *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, May 2015.



- [27] P. P. S. Narayan: “Sherpa Update,” *developer.yahoo.com*, June 8.
- [28] Martin Kleppmann: “Bottled Water: Real-Time Integration of PostgreSQL and Kafka,” *martin.kleppmann.com*, April 23, 2015.
- [29] Ben Osheroﬀ: “Introducing Maxwell, a mysql-to-kafka Binlog Processor,” *developer.zendesk.com*, August 20, 2015.
- [30] Randall Hauch: “Debezium 0.2.1 Released,” *debezium.io*, June 10, 2016.
- [31] Prem Santosh Udaya Shankar: “Streaming MySQL Tables in Real-Time to Kafka,” *engineeringblog.yelp.com*, August 1, 2016.
- [32] “Mongoriver,” Stripe, Inc., *github.com*, September 2014.
- [33] Dan Harvey: “Change Data Capture with Mongo + Kafka,” at *Hadoop Users Group UK*, August 2015.
- [34] “Oracle GoldenGate 12c: Real-Time Access to Real-Time Information,” Oracle White Paper, March 2015.
- [35] “Oracle GoldenGate Fundamentals: How Oracle GoldenGate Works,” Oracle Corporation, *youtube.com*, November 2012.
- [36] Slava Akhmechet: “Advancing the Realtime Web,” *rethinkdb.com*, January 27, 2015.
- [37] “Firebase Realtime Database Documentation,” Google, Inc., *firebase.google.com*, May 2016.
- [38] “Apache CouchDB 1.6 Documentation,” *docs.couchdb.org*, 2014.
- [39] Matt DeBergalis: “Meteor 0.7.0: Scalable Database Queries Using MongoDB Oplog Instead of Poll-and-Diff,” *info.meteor.com*, December 17, 2013.
- [40] “Chapter 15. Importing and Exporting Live Data,” VoltDB 6.4 User Manual, *docs.voltdb.com*, June 2016.
- [41] Neha Narkhede: “Announcing Kafka Connect: Building Large-Scale Low-Latency Data Pipelines,” *confluent.io*, February 18, 2016.

- [42] Greg Young: “CQRS and Event Sourcing,” at *Code on the Beach*, August 2014.
- [43] Martin Fowler: “Event Sourcing,” *martinfowler.com*, December 12, 2005.
- [44] Vaughn Vernon: *Implementing Domain-Driven Design*. Addison-Wesley Professional, 2013. ISBN: 978-0-321-83457-7.
- [45] H. V. Jagadish, Inderpal Singh Mumick, and Abraham Silberschatz: “View Maintenance Issues for the Chronicle Data Model,” at *14th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, May 1995. doi:10.1145/212433.220201.
- [46] “Event Store 3.5.0 Documentation,” Event Store LLP, *docs.geteventstore.com*, February 2016.
- [47] Martin Kleppmann: *Making Sense of Stream Processing*. Report, O’ Reilly Media, May 2016.
- [48] Sander Mak: “Event-Sourced Architectures with Akka,” at *JavaOne*, September 2014.
- [49] Julian Hyde: personal communication, June 2016.
- [50] Ashish Gupta and Inderpal Singh Mumick: *Materialized Views: Techniques, Implementations, and Applications*. MIT Press, 1999. ISBN: 978-0-262-57122-7.
- [51] Timothy Griffin and Leonid Libkin: “Incremental Maintenance of Views with Duplicates,” at *ACM International Conference on Management of Data (SIGMOD)*, May 1995. doi:10.1145/223784.223849.
- [52] Pat Helland: “Immutability Changes Everything,” at *7th Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2015.
- [53] Martin Kleppmann: “Accounting for Computer Scientists,” *martin.kleppmann.com*, March 7, 2011.
- [54] Pat Helland: “Accountants Don’t Use Erasers,” *blogs.msdn.com*, June 14, 2007.
- [55] Fangjin Yang: “Dogfooding with Druid, Samza, and Kafka: Metametrics at Metamarkets,” *metamarkets.com*, June 3, 2015.



- [56] Gavin Li, Jianqiu Lv, and Hang Qi: “Pistachio: Co-Locate the Data and Compute for Fastest Cloud Compute,” *yahoohadoop.tumblr.com*, April 13, 2015.
- [57] Kartik Paramasivam: “Stream Processing Hard Problems - Part 1: Killing Lambda,” *engineering.linkedin.com*, June 27, 2016.
- [58] Martin Fowler: “CQRS,” *martinfowler.com*, July 14, 2011.
- [59] Greg Young: “CQRS Documents,” *cqrs.files.wordpress.com*, November 2010.
- [60] Baron Schwartz: “Immutability, MVCC, and Garbage Collection,” *xaprb.com*, December 28, 2013.
- [61] Daniel Eloff, Slava Akhmechet, Jay Kreps, et al.: “Re: Turning the Database Inside-out with Apache Samza,” Hacker News discussion, *news.ycombinator.com*, March 4, 2015.
- [62] “Atomic Development Resources: Excision,” Cognitect, Inc., *docs.datomic.com*.
- [63] “Fossil Documentation: Deleting Content from Fossil,” *fossil-scm.org*, 2016.
- [64] Jay Kreps: “The irony of distributed systems is that data loss is really easy but deleting data is surprisingly hard,” *twitter.com*, March 30, 2015.
- [65] David C. Luckham: “What’s the Difference Between ESP and CEP?,” *complexevents.com*, August 1, 2006.
- [66] Srinath Perera: “How Is Stream Processing and Complex Event Processing (CEP) Different?,” *quora.com*, December 3, 2015.
- [67] Arvind Arasu, Shivnath Babu, and Jennifer Widom: “The CQL Continuous Query Language: Semantic Foundations and Query Execution,” *The VLDB Journal*, volume 15, number 2, pages 121-142, June 2006. doi:10.1007/s00778-004-0147-z.
- [68] Julian Hyde: “Data in Flight: How Streaming SQL Technology Can Help Solve the Web 2.0 Data Crunch,” *ACM Queue*, volume 7, number 11, December 2009. doi:10.1145/1661785.1667562.
- [69] “Esper Reference, Version 5.4.0,” EsperTech, Inc., *espertech.com*, April 2016.





- [70] Zubair Nabi, Eric Bouillet, Andrew Bainbridge, and Chris Thomas: “Of Streams and Storms,” IBM technical report, *developer.ibm.com*, April 2014.
- [71] Milinda Pathirage, Julian Hyde, Yi Pan, and Beth Plale: “SamzaSQL: Scalable Fast Data Management with Streaming SQL,” at *IEEE International Workshop on High-Performance Big Data Computing (HPBDC)*, May 2016. doi:10.1109/IPDPSW.2016.141.
- [72] Philippe Flajolet, Éric Fusy, Olivier Gandouet, and Frédéric Meunier: “HyperLog Log: The Analysis of a Near-Optimal Cardinality Estimation Algorithm,” at *Conference on Analysis of Algorithms (AofA)*, June 2007.
- [73] Jay Kreps: “Questioning the Lambda Architecture,” *oreilly.com*, July 2, 2014.
- [74] Ian Hellström: “An Overview of Apache Streaming Technologies,” *databaseline.wordpress.com*, March 12, 2016.
- [75] Jay Kreps: “Why Local State Is a Fundamental Primitive in Stream Processing,” *oreilly.com*, July 31, 2014.
- [76] Shay Banon: “Percolator,” *elastic.co*, February 8, 2011.
- [77] Alan Woodward and Martin Kleppmann: “Real-Time Full-Text Search with Luwak and Samza,” *martin.kleppmann.com*, April 13, 2015.
- [78] “Apache Storm 1.0.1 Documentation,” *storm.apache.org*, May 2016.
- [79] Tyler Akidau: “The World Beyond Batch: Streaming 102,” *oreilly.com*, January 20, 2016.
- [80] Stephan Ewen: “Streaming Analytics with Apache Flink,” at *Kafka Summit*, April 2016.
- [81] Tyler Akidau, Alex Balikov, Kaya Bekiroğlu, et al.: “MillWheel: Fault-Tolerant Stream Processing at Internet Scale,” at *39th International Conference on Very Large Data Bases (VLDB)*, August 2013.
- [82] Alex Dean: “Improving Snowplow’s Understanding of Time,” *snowplowanalytics.com*, September 15, 2015.
- [83] “Windowing (Azure Stream Analytics),” Microsoft Azure Reference, *msdn.microsoft.com*, April 2016.



- [84] “State Management,” Apache Samza 0.10 Documentation, *samza.apache.org*, December 2015.
- [85] Rajagopal Ananthanarayanan, Venkatesh Basker, Sumit Das, et al.: “Photon: Fault-Tolerant and Scalable Joining of Continuous Data Streams,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2013. doi: 10.1145/2463676.2465272.
- [86] Martin Kleppmann: “Samza Newsfeed Demo,” *github.com*, September 2014.
- [87] Ben Kirwin: “Doing the Impossible: Exactly-Once Messaging Patterns in Kafka,” *ben.kirw.in*, November 28, 2014.
- [88] Pat Helland: “Data on the Outside Versus Data on the Inside,” at *2nd Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2005.
- [89] Ralph Kimball and Margy Ross: *The Data Warehouse Toolkit: The Definitive Guide to Dimensional Modeling*, 3rd edition. John Wiley & Sons, 2013. ISBN: 978-1-118-53080-1.
- [90] Viktor Klang: “I’m coining the phrase ‘effectively-once’ for message processing with at-least-once + idempotent operations,” *twitter.com*, October 20, 2016.
- [91] Matei Zaharia, Tathagata Das, Haoyuan Li, et al.: “Discretized Streams: An Efficient and Fault-Tolerant Model for Stream Processing on Large Clusters,” at *4th USENIX Conference in Hot Topics in Cloud Computing (HotCloud)*, June 2012.
- [92] Kostas Tzoumas, Stephan Ewen, and Robert Metzger: “High-Throughput, Low-Latency, and Exactly-Once Stream Processing with Apache Flink,” *data-artisans.com*, August 5, 2015.
- [93] Paris Carbone, Gyula Fóra, Stephan Ewen, et al.: “Lightweight Asynchronous Snapshots for Distributed Dataflows,” *arXiv:1506.08603 [cs.DC]*, June 29, 2015.
- [94] Ryan Betts and John Hugg: *Fast Data: Smart and at Scale*. Report, O’Reilly Media, October 2015.
- [95] Flavio Junqueira: “Making Sense of Exactly-Once Semantics,” at *Strata+Hadoop World London*, June 2016.

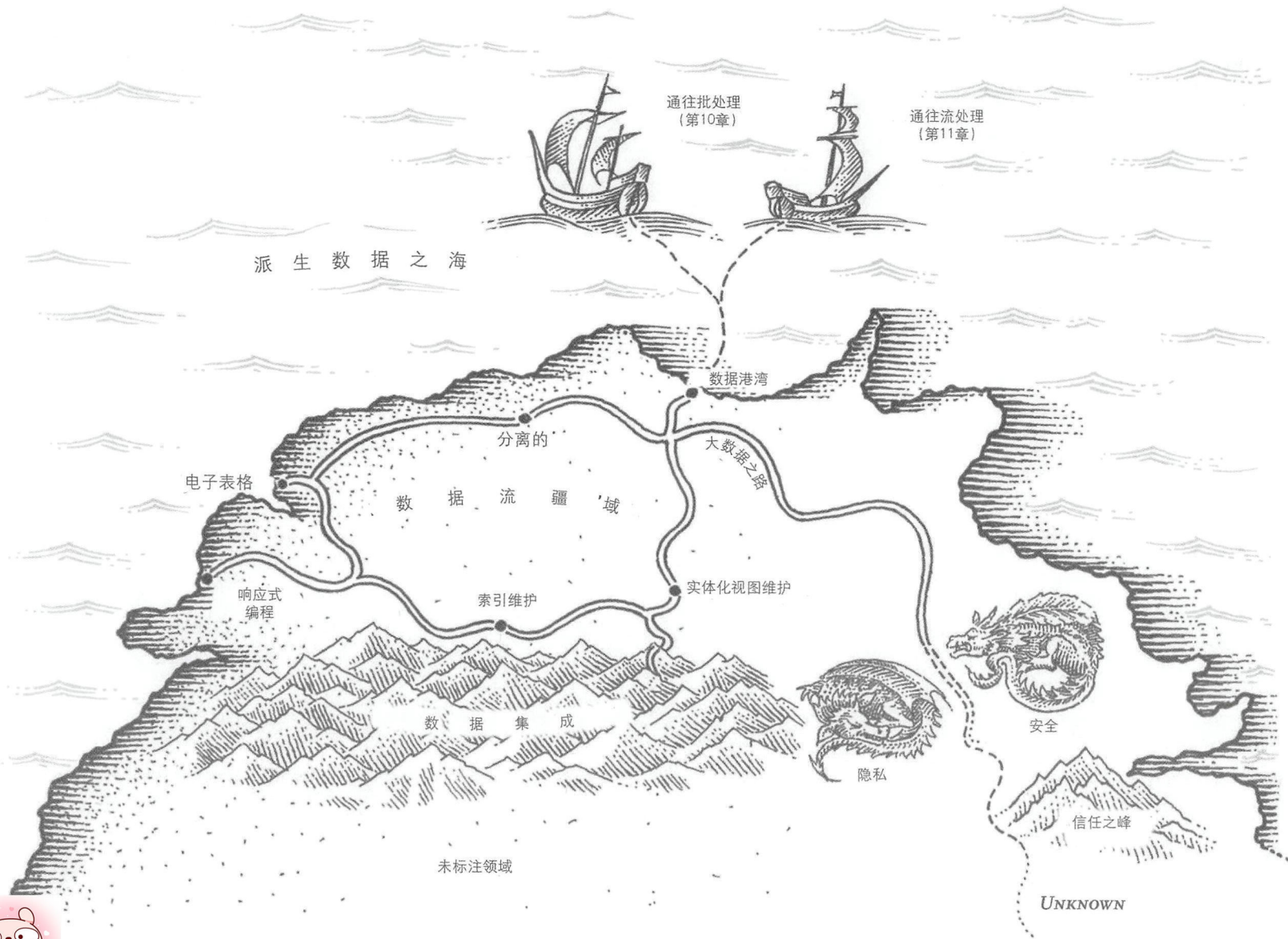




- [96] Jason Gustafson, Flavio Junqueira, Apurva Mehta, Sriram Subramanian, and Guozhang Wang: “KIP-98-Exactly Once Delivery and Transactional Messaging,” *cwiki.apache.org*, November 2016.
- [97] Pat Helland: “Idempotence Is Not a Medical Condition,” *Communications of the ACM*, volume 55, number 5, page 56, May 2012. doi:10.1145/2160718.2160734.
- [98] Jay Kreps: “Re: Trying to Achieve Deterministic Behavior on Recovery/Rewind,” email to *samza-dev* mailing list, September 9, 2014.
- [99] E. N. (Mootaz) Elnozahy, Lorenzo Alvisi, Yi-Min Wang, and David B. Johnson: “A Survey of Rollback-Recovery Protocols in Message-Passing Systems,” *ACM Computing Surveys*, volume 34, number 3, pages 375-408, September 2002. doi: 10.1145/568522.568525
- [100] Adam Warski: “Kafka Streams - How Does It Fit the Stream Processing Landscape?,” *softwaremill.com*, June 1, 2016.









# 数据系统的未来

如果船长的最高目标是保住他的船，那么他只能永远待在港口。

——St. Thomas Aquinas，神学大全（1265-1274）

到目前为止，这本书主要描述了当前流行的技术。在最后一章中，我们将谈谈对未来的看法，未来的系统应该是什么样子的。我将提出一些想法和方法，我相信这些方法可能会从根本上改进我们设计和构建应用程序的方式。

当然，这些意见和猜测带有主观性，所以本章当阐述的是我个人观点时会使用第一人称。如果与你的看法和理念有所不同，非常欢迎提出不同的意见。但我希望本章的观点至少可以成为一个富有成效的讨论的出发点，并澄清一些经常混淆的概念。

开篇第1章总结了本书的目标：探索如何构建可靠、可扩展和可维护的应用和系统。这些主题贯穿了之后的所有章节。例如，我们讨论了许多有助于提高可靠性的容错算法，提高可扩展性的分区，以及为提高可维护性相关演化与抽象机制。本章，我们将把所有这些想法放在一起，并以此为基础来展望未来。我们的目标是探索如何设计比今天更好的应用程序：强大、正确、可演化并最终使所有人受益。

## 数据集成

本书反复强调了一个主题，对于任何给定的问题，都有多种解决方案，而这些解决方案都有各自优缺点和折中之处。例如，在第3章讨论存储引擎时，我们看到了日志结构存储，B-tree和列式存储。在第5章讨论复制时，我们介绍了主从复制，多主节点复制和无主节点复制。





如果给你一个问题，比如“需要把数据保存起来，可以支持查询”，这样的问题其实不存在唯一正确的解决方案，或者说不同具体情况下有不同的方案。而软件的实现通常必须选择一种特定的方法。通常一种实现很难同时兼具鲁棒与性能，如果试图兼顾所有方面，那唯一可以肯定的是最终的实现一定很糟糕。

因此，选择合适的软件组件也需要视情况而定。每一个软件，即使是所谓的“通用”数据库，也都是针对特定的使用模式而设计的。

面对如此众多的备选方案，第一个挑战就是弄清楚软件产品与他们适合运行环境之间的对应关系。软件厂商是不愿意告诉你它们的软件不适合哪些负载，所以系统通过前面的章节可以帮助你准备一些有针对性的问题，例如产品字里行间的含义以及设计取舍之道。

然而，即使你完全理解了工具与使用环境之间的对应关系，还有一个挑战摆在那里：在复杂的应用程序中，数据通常以多种不同的方式被使用。不太可能存在适用于所有不同环境的软件，因此你不可避免地要将几个不同的软件组合在一起，以提供应用程序的功能性。

## 采用派生数据来组合工具

例如，经常需要将一个OLTP数据库和全文搜索索引集成起来以处理任意关键字查询。虽然有些数据库（如PostgreSQL）支持全文索引功能，可以满足一些简单应用<sup>[1]</sup>，但更复杂的查询则需要专业的信息检索工具。另一方面，搜索索引通常不太适合作为一个持久的记录系统，所以许多应用程序需要结合两种不同的工具来满足所有需求。

我们在第11章“保持系统同步”中提到了数据系统集成问题。随着不同类型的数据持续增加，集成问题会变得越来越困难。除了数据库和搜索索引，你可能需要在分析系统中保存数据副本（例如数据仓库，批处理和流处理系统）；从派生数据中维护数据缓存或非规范化版本；把数据传递给机器学习系统，或者分类，排序以及推荐系统；还有发送数据更改的通知等。

令人惊讶的是，我经常看到软件工程师煞有介事地说，“根据我的经验，99%的人只需要X”或“不需要X”（X可以代表各种值）。我认为这样的说法更多地在他多么多么的有经验，而不是在谈论某项技术的实际用途。针对数据所做的事情千变万化，涉及的范围也是惊人的广泛。某些人认为是一个模糊而毫无意义的特性对其他人来说可能是核心需求。数据集成的需求通常只有在缩小并考虑整个组织框架内数据流时才





会变得更加凸显。

## 为何需要数据流

当同样数据的多个副本需要保存在不同的存储系统，以满足不同的访问模式需求时，你需要搞清楚数据输入和输出：数据第一次是写在哪里？不同的数据拷贝来自哪些数据源？如何以正确的格式将数据导入到所有正确的位置？

我们举个例子，比如可以将数据首先写入记录系统数据库，捕获对该数据库的更改（参见第11章“变更数据捕获”），然后按相同顺序将更改应用到搜索索引中。如果变更数据捕获（CDC）是更新索引的唯一方法，那么你就可以确信索引完全来自记录系统，因此是一致的（软件bug除外）。写数据库是向系统提供新输入的唯一途径。

允许应用程序同时向搜索索引和数据库写入会带来问题（见图11-4），两个客户端同时发送冲突的写操作，而两个存储系统以不同的顺序来执行它们。在这种情况下，负责确定写操作顺序的既不是数据库也不是搜索索引，因此它们可能会做出相互矛盾的决定，以至于导致彼此之间永久的不一致。

如果可以通过单个系统来决定所有输入的写入顺序，那么以相同的顺序处理写操作就可以更容易地派生出数据的其他表示形式。这就是第9章“全序关系广播”中所介绍的状态机复制方法。无论是使用变更数据捕获还是事件获取日志，都不如简化总体顺序的原则重要。

根据事件日志来更新一个派生数据系统通常会比较好实现，并且可以实现确定性和幂等性（幂等概念参见第11章），也因此使系统很容易从故障中恢复。

## 派生数据与分布式事务

为了保持不同的数据系统彼此之间的一致性，经典的方法是通过分布式事务，我们在第9章讨论过原子提交和两阶段提交机制（2PC）。那么与分布式事务相比，派生数据系统的方法怎么样呢？

抽象点说，它们通过不同的方式达到类似的目标。分布式事务通过使用锁机制进行互斥来决定写操作的顺序（请参阅第7章“两阶段加锁”），而CDC和事件源使用日志进行排序。分布式事务使用原子提交来确保更改只生效一次，而基于日志的系统通常基于确定性重试和幂等性。

最大的不同在于事务系统通常提供线性化（请参阅第9章“线性化”），这意味着它可以保证读自己的写等一致性（请参阅第5章“读自己的写”）。另一方面，派生的



数据系统通常是异步更新的，所以默认情况下它们无法提供类似级别保证。

在能够承受付出分布式事务所带来开销的特定环境中，分布式事物已有成功案例。但是，我认为XA的容错性和性能不尽如人意（请参阅第9章的“实践中的分布式事务”），这严重限制了它的实用性。我相信为分布式事务未来可以有更好的协议，但是这样的协议被广泛采用而且能与现有工具集成将非常具有挑战性，并且不太可能在短期内发生。

考虑到好的分布式事务协议尚未得到广泛支持，我认为基于日志的派生数据是集成不同数据系统的最有前途的方法。然而，一些保证如“读自己写”仍是非常有用，因此我不认为到处宣扬“最终的一致性是不可能的，接受并学会处理它”有任何实际意义（至少在没有很好的指导来告诉人们如何处理它的情况下是不太可行的）。

在本章后面“正确性目标”中，我们将讨论在异步派生系统之上实现更强保证的一些方法，并在分布式事务和基于日志的异步系统之外找到一条中间的路。

## 全序的局限

对于非常小的系统，构建一个完全有序的事件日志是完全可行的（流行的主从复制正是构造了这样的日志）。但是，随着系统越来越大，并且面对更为复杂的负载时，瓶颈就开始出现了：

- 在大多数情况下，构建一个完全有序的日志需要所有事件都通过一个主节点来决定排序。如果事件吞吐量大于单台节点可处理的上限，则需要将其分区到多台节点上（请参见第11章的“分区日志”），这就使得两个不同分区中的事件顺序变得不明确了。
- 如果服务器分布在多个不同地理位置的数据中心，为了避免整个数据中心不可用，且考虑到网络延迟使跨数据中心协调的同步效率很低，因此通常在每个数据中心都有独立的主节点（参阅第5章“多主节点复制”）。这意味着来自两个不同数据中心的事件顺序不确定。
- 将应用程序部署为微服务时（请参阅第4章“基于服务的数据流：REST和RPC”），常见的设计是将每个服务与其持久化的状态一起作为独立单元部署，而服务之间不共享持久化状态。当两个事件来自不同的服务时，这些事件没有清楚的顺序。
- 某些应用程序在客户端维护一些状态，当用户输入时会立即更新（无需等待服务器的确认），甚至可以继续离线工作（参阅第5章“客户端的离线操作”）。对于这样的应用程序，客户端和服务器很可能看到不同的事件顺序。



从形式上讲，决定事件的全序关系称为全序关系广播，它等价于共识（参见第9章“共识算法与全序关系广播”）。大多数共识算法是针对单节点吞吐量足以处理整个事件流而设计的，并且这些算法不提供支持多节点共享事件排序的机制。设计突破单节点吞吐量甚至在广域地理环境分布的共识算法仍然是一个有待研究的开放性问题。

## 排序事件以捕获因果关系

如果事件之间不存在因果关系，则不支持全序排序并不是一个大问题，因为并发事件可以任意排序。其他一些情况很容易处理：例如，当同一对象有多个更新时，可以通过将特定对象ID的所有更新都路由到同一个日志分区上。然而，因果依赖有时候会有一些更微妙的情况（详见第9章“排序与因果关系”）。

例如在一个社交网络服务里，有两个以前相恋，但现在分手了的用户。女孩子将这个分手的男朋友删除后，发送消息给她其余的朋友抱怨她的前男友。这个女孩子以为，他的前男友不应该看到这些抱怨信息，因为该消息是她删除男朋友的朋友状态之后发送的。

然而，如果朋友状态和消息分别存储在不同地方，那么可能会丢失“删除好友”事件和“抱怨”事件之间的顺序关系。如果没有正确捕获事件的因果关系，负责发送新消息通知的服务可能在“删除好友”之前就把“抱怨”消息广播给所有朋友，结果就是其前恋人还是会收到一条本不应该看到的消息。

在这个例子中，通知实际上是消息和朋友列表的联结操作，这和之前讨论的join与时间问题有关系了（请参阅第11章“join的时间依赖性”）。不幸的是，似乎这个问题并没有一个简单的答案<sup>[2, 3]</sup>。从头分析的话是这样的：

- 逻辑时间戳可以在无协调者情况下提供的全序关系（参见第9章“序列号排序”），所以当全序关系广播不可行时可以用得上。但是，它们仍然需要接收者去处理那些乱序事件，并且需要额外的元数据。
- 如果可以记录一条事件来标记用户在做决定以前所看到系统状态，并给该事件一个唯一的标识符，那么任何后续的事件都可以通过引用该事件标识符来记录因果关系<sup>[4]</sup>。我们将在本章后面“读操作也是事件”中重谈这个概念。
- 冲突解决算法（请参阅第5章的“自动冲突解决”）可以处理异常顺序的事件。他们可以帮助维护正确状态，但前提是行为不能外部可见的副作用（例如向用户发送通知）。



也许将来应用程序开发的模式将会出现可喜变化，使得因果关系被有效地捕获，派生状态得到正确维护，而不会强制所有事件必须通过全序广播（系统瓶颈）。

## 批处理和流处理集成

我认为数据整合的目标是确保数据在所有正确的地方以正确的形式结束。这样做涉及消费输入数据，转换，join，过滤，聚合，训练模型，评估并最终写入适当的输出。而批处理和流处理则是实现这一目标的有效工具。

批处理和流处理的输出是派生数据集，如搜索索引，实体化视图，向用户展现的建议，聚合度量标准等（请参阅第10章“批处理工作流的输出”和第11章的“流处理的适用场景”）。

正如我们在第10章和第11章中所介绍的，批处理和流处理有许多共同的原则，而根本区别在于流处理器运行在无界数据集上，而批处理的输入是已知的有限大小。处理引擎的实现方式也有很多细节上的差异，但是这些区别现在变得越来越模糊。

Spark通过将流分解为微批处理来在批处理引擎之上执行流处理，而Apache Flink则直接在流处理引擎上执行批处理<sup>[5]</sup>。原则上，一种类型的处理可以通过另一种类型来模拟，尽管性能特征会有所不同：例如，微批处理可能在跳跃或滑动窗口上表现不佳<sup>[6]</sup>。

### 保持派生状态

批处理具有相当强的功能特性（即使代码不是用函数式编程语言编写的），包括倡导确定性、纯函数操作即输出仅依赖于输入，除了显式输出以外没有任何副作用，输入不可变，追加式输出结果等。流处理是类似的，但它扩展了操作来支持可管理的、容错的状态（请参阅第11章“故障后重建状态”）。

具有良好定义的输入和输出的确定性函数原理上不仅有利于容错（请参阅第11章“幂等性”），而且还简化了组织中数据流的推理<sup>[7]</sup>。无论派生数据是搜索索引，统计模型还是缓存，从数据管道的角度来看，对于从一个事物派生出另一个事物，通过功能应用程序代码推动一个系统中的状态更改以及将这种效果应用到派生系统，都是有帮助的。

原则上，派生数据系统可以同步维护，就像关系型数据库在同一个事务中同步地更新二级索引，就像写入被索引的表一样。然而，异步是使基于事件日志的系统健壮的原因：它允许系统的一部分在本地包含故障，而如果任何参与者失败，则分布式事务中

止，因此他们倾向于通过将故障扩展到其余部分来放大故障（请参阅第9章“分布式事务的限制”）。

我们在第6章“分区与二级索引”中看到，二级索引经常跨越分区边界。具有二级索引的分区系统需要将写入发送到多个分区（如果索引文件分区）或将读取发送到所有分区（如果索引是文档分区的话）。如果索引是异步维护的，这种交叉分区通信也是最可靠的和最易扩展的<sup>[8]</sup>（另请参阅本章后面“多分区数据处理”）。

## 为应用程序演化而重新处理数据

在需要维护派生数据时，批处理和流处理都会用得上。流处理可以将输入的变化数据迅速反映在派生视图中，而批处理则可以反复处理大量的累积数据，以便将新视图导出到现有数据集上。

特别是对现有数据进行重新处理，为维护系统提供了一个良好的机制，平滑支持新功能以及多变的需求（参见第4章）。如果不进行重新处理，则只能支持有限简单的模式变化，比如将新的可选字段添加到记录中，或者添加新的记录类型。这些简单的模式比较常见与写时模式和读时模式（参阅第2章“文档模型中的模式灵活性”）。另一方面，通过重新处理，可以将数据集重组为一个完全不同的模型，以便更好地满足新要求。

### 发生在铁路上的模式迁移

大规模的“模式迁移”并发仅限于计算机系统。例如，19世纪英国铁路建设的初期，轨距（两轨之间的距离）存在各种竞争标准。按照某一个尺寸标准而建的列车无法在另一个尺寸标准的轨道上运行，这严重限制了铁路网络大规模互连<sup>[9]</sup>。

1846年最终确定了一个国家标准，其他标准的轨道因此必须更换，但是如何在关闭火车线路的情况下进行数月甚至数年的调整呢？解决的办法是通过增加第三条铁轨的方法先把轨道换成双轨或混合轨距。这种转换可以逐步完成，当完成后，新老标准的列车都可以选择三条导轨的两条在线上运行。最后，一旦所有列车都支持标准轨距，则可以拆除提供非标准轨距的轨道。

以这种方式“再加工”现有的轨道，让新旧版本并存，可以在几年的时间内逐步改变轨距。然而，这是一个代价昂贵的项目，这就是今天仍然看到很多非标准轨道原因。例如，旧金山湾区的BART系统仍在使用与美国大部分地区不同的标准。

派生视图允许逐步演变。如果想重新构建数据集，无需采用高风险的陡然切换。而是可以在同一个基础数据上的两个独立派生视图来同时维护新老两种架构。然后，逐步开始将少量用户迁移到新视图中，以测试其性能并发现是否有错误，而大多数用户将继续路由到旧视图。之后，逐渐增加访问新视图的用户比例，最终放弃旧视图<sup>[10]</sup>。

这种逐渐迁移的美妙之处在于，如果出现问题，每个阶段的过程都是可以轻易反转：你总是有一个工作系统可以回退。通过减少不可逆损害的风险，可以更加自信的向前推进，从而更快地改善系统<sup>[11]</sup>。

## Lambda架构

如果使用批处理来重新处理历史数据，并且使用流处理来处理最近的更新，那么如何将这两者结合起来呢？备受关注的Lambda结构<sup>[12]</sup>是一个不错的解决方案。

Lambda体系结构的核心思想是进来的数据以不可变事件形式追加写到不断增长的数据集，类似于事件源（请参阅第11章的“事件溯源”）。基于这些总事件，可以派生出读优化的视图。lambda结构建议并行运行两个不同的系统：一个批处理系统如Hadoop MapReduce，以及一个单独的流处理系统，如Storm。

在lambda方法中，流处理器处理事件并快速产生对试图的近似更新；批处理系统则处理同一组事件并产生派生视图的校正版本。这种设计背后的原因是，批处理更简单，不易出现错误，而流处理器则相对太可靠的，难以实现容错（请参阅第11章“流处理的容错”）。而且流处理可以使用快速的近似算法，而批处理只能用较慢的精确算法。

Lambda架构师一个很有影响力的想法，它将数据系统的设计变得更好，特别是推广了将派生视图用于不可变事件流以及必要时重新处理事件的原则。但是，另一方面，我认为它有一些实际问题：

- 不得不在批处理和流处理框架中运行相同的处理逻辑，这是一项重大的额外工作。虽然像Summingbird这样的库<sup>[13]</sup>提供了一个可以在批处理或流处理环境下运行的计算抽象，但调试/优化和维护两个不同系统的操作复杂性依然存在<sup>[14]</sup>。
- 由于流处理和批处理各自产生单独的输出，因此需要合并二者结果再响应给用户请求。如果计算是通过滚动窗口的简单聚合，则此合并相当容易，但如果包含更复杂的操作（例如连接和会话流程）导出视图，或者输出不是时间序列，则合并非常困难。
- 能够重新处理整个历史数据集思路是好，但是在大型数据集上这样做往往代价太



高。因此，批处理流水线通常需要设置为处理增量批处理（例如，在每小时结束时处理这一小时的数据），而不是重新处理所有事情。这会引发第11章“流的时间问题”中处理滞后消息和处理跨批量边界的窗口等问题。增加批量计算就增加了复杂性，使其更类似于流式，这与保持批处理层尽可能简单的目标背道而驰。

## 统一批处理和流处理

最近的发展使得lambda结构可以充分发挥其优点而规避其缺点，那就是允许批处理计算（重新处理历史数据）和流计算（事件到达时即处理）在同一个系统中实现<sup>[15]</sup>。

在一个系统中统一批处理和流处理需要以下功能，而这些功能目前正越来越普及：

- 支持以相同的处理引擎来处理最新事件和处理历史回放事件。例如，基于日志的消息代理可以重放消息（请参阅第11章“重新处理消息”），某些流处理器可以从诸如HDFS等分布式文件系统读取输入。
- 支持只处理一次语义。即使事实上发生了错误（请参阅第11章“流处理的容错”），流处理器依然确保最终的输出与未发生错误的输出相同。与批处理一样，这需要丢弃任何失败任务的部分输出。
- 支持依据事件发生时间而不是处理时间进行窗口化。因为在重新处理历史事件时，处理时间毫无意义（请参阅第11章“流的时间问题”）。例如，Apache Beam提供了用于表示这种计算的API，支持在Apache Flink或Google Cloud Dataflow引擎上运行。

## 分拆数据库

从最抽象的层面上理解，数据库、Hadoop和操作系统都提供了相同的功能：它们保存某些数据，并支持处理和查询数据<sup>[16]</sup>。数据库将数据存储在某些数据模型的记录中（表中的行，文档，图中的点等），而操作系统的文件系统将数据存储的文件中，但是核心都是“信息管理”系统<sup>[17]</sup>。正如我们在第10章中介绍的，Hadoop生态系统有点像UNIX的分布式版本。

当然，有存在很大实际的区别。例如，许多文件系统不能很好地处理包含一千万个小文件的目录，而包含一千万条记录的数据库是完全正常的。尽管如此，操作系统和数据库之间的异同之处仍是值得探讨的。

UNIX和关系型数据库采用了大不一样的哲学思想看待信息管理问题。UNIX认为它的目的是为程序员提供一个逻辑的，但是相当低层次的硬件抽象，而关系型数据库则



希望为应用程序员提供一个高层次的抽象，来隐藏磁盘上数据结构的复杂性、并发性、崩溃恢复等。UNIX开发的管道和文件只是字节序列，而数据库开发了SQL和事务。

哪种方法更好？当然，这取决于你想要什么。UNIX是“简单的”，因为它是硬件资源的一层包装；关系型数据库也是“简单的”，因为一个很短的声明性查询可以利用很多强大的基础设施（查询优化、索引、join方法、并发控制、复制等），而无需发起请求者详细理解实现细节。

这些哲学思想之间的对立已经持续了几十年（UNIX和关系模型都出现在20世纪70年代初），目前仍然没有解决。例如，我宁愿将NoSQL的发展解释为一种将低级别抽象方法应用于分布式OLTP领域的诉求。

在这一部分，我试图辩证看待两种哲学，希望很好地将二者结合在一起，做到两全其美。

## 编排多种数据存储技术

在本书中，我们讨论了数据库提供的各种功能及其工作原理，其中包括：

- 二级索引：根据字段值高效地搜索所有记录（请参阅第3章“其他索引结构”）。
- 实体化视图：预先计算查询结果并将其缓存（请参阅第3章“聚合：数据立方体与物化视图”）。
- 复制日志：使多节点上数据副本保持最新（请参阅第5章“复制日志的实现”）。
- 全文搜索索引：在文本中进行关键字搜索（请参阅第3章“全文搜索与模糊索引”）并且内置于某些关系型数据库<sup>[1]</sup>。

在第10章和第11章中介绍了类似的主题。我们讨论了如何构建全文搜索索引（请参阅第10章“批处理工作流的输出”），实体化视图维护（请参阅第11章“维护实体化视图”），以及如何将变更数据从数据库复制到派生数据系统（请参阅第11章的“变更数据捕获”）。

由此看来，数据库中内置的这些功能与采用批处理和流处理器构建的派生数据系统似乎很多对应关系。



## 创建一个索引

想想在关系型数据库中运行CREATE INDEX来创建新索引时会发生什么。数据库必须扫描表的一致性快照，挑选出所有被索引的字段值，对它们进行排序，然后得到索引。接下来，必须处理从一致性快照创建以来所累计的写入操作（假设表在创建索引时未被锁定，所以写操作可能会继续）。完成后，只要有事务写入表中，数据库就必须持续保持索引处于最新状态。

这个过程和配置新的从节点副本非常相似（请参阅第5章的“配置新的从节点”），也非常类似于流处理系统中的初始变更数据捕获（请参阅第11章“初始化快照”）。

无论何时运行CREATE INDEX，数据库都会重新处理现有的数据集（正如我们在本章前面“为应用程序演化而重新处理数据”中所述），并将索引作为新视图导出到现有数据上。现有数据可能是状态的快照，而不是所有发生变化的日志，但两者却密切相关（请参阅第11章“状态，流与不可变性”）。

## 元数据库

有鉴于此，我认为整个组织的数据流开始变得像一个巨大的数据库<sup>[7]</sup>。每当批处理，流或ETL过程将数据从某个一个位置（和表单）传输到另一个位置（或表单）时，它就像数据库子系统一样需要保持索引或实体化视图至最新状态。

这样看来，批处理和流处理器就像触发器，存储过程和实体化视图维护相关实现。它们所维护的派生数据系统就像不同的索引类型。例如，关系型数据库可能支持B-Tree索引、哈希索引、空间索引（请参阅第3章“多列索引”）以及其他类型的索引。在新兴的派生数据系统体系结构中，这些特性不是靠单一集成数据产品来实现，而是由各种不同的软件所提供，他们运行在不同的机器上，并由不同的团队来管理。

这些发展在未来会我们引向何处？如果我们的前提是，没有一个统一的数据模型或存储格式适用于所有的访问模式，我推测，到时候会有两种途径，不同的存储和处理工具最终会组合成一个紧密结合的系统：

### 联合数据库：统一读端

可以为各种各样的底层存储引擎和处理方法提供一个统一的查询接口：一种称为联合数据库或聚合存储的方法<sup>[18,19]</sup>。例如，PostgreSQL的外部数据包功能就符合这种模式<sup>[20]</sup>。需要专门的数据模型或查询接口的应用程序仍然可以直接访问底层存储引擎，而想要组合来自不同位置的数据的用户可以通过联合接口轻松完成任务。







联合查询接口遵循单一集成系统的关系模型，并具有高级查询语言和优雅语义，但是内部实现非常复杂。

### 分离式数据库：统一写端

虽然上述联合式方案能解决跨系统的只读查询问题，但是跨系统的同步写入方面却没有一个好的答案。我们说对于一个数据库，创建一致的索引是其内置的功能。而在构建跨多个存储系统的数据库时，我们同样需要确保所有数据更改都会体现在所有正确的位置上，即使中间发生了某些故障。多个存储系统可以可靠地连接在一起（例如，通过变更数据捕获和事件日志），类似思路，可以跨不同技术方案来同步写，从而实现分离式（或非捆绑式）数据库的索引维护功能<sup>[7,21]</sup>。

这种分离的方法正是遵循了UNIX传统：单一任务做好一件事，且内部通过统一的低级API（管道）进行通信，外部通过更高级别的语言（shell）<sup>[16]</sup>进行组合。

## 分离式如何工作

联合方式与分离方式可以看出同一个硬币的两面：用不同的组件构成一个可靠、可扩展的和可维护的系统。联合的只读查询需要将一个数据模型映射到另一个数据模型，这里需要对模型匹配方面足够深入的考虑，但最终它还是一个可控问题。我认为同步写入多个存储系统是一个更难的工程问题，所以我将重点关注这个问题。

传统的同步写依赖于跨异构存储系统的分布式事务<sup>[18]</sup>，我认为这是个错误的解决方案（请参阅本章前面“派生数据与分布式事务”）。在单个存储系统内或流处理系统内的事务是可行的，但是当数据跨越不同技术的边界时，我认为具有幂等写入的异步事件日志是一种更加健壮和可行的方法。

例如，在一些流处理器中使用分布式事务来实现exactly-once语义（请参阅第11章“重新思考原子提交”），这是可行的。但是，当一个事务涉及到多个不同人群负责的系统时（例如，当数据从流系统写入到分布式键值存储或搜索索引时），缺乏标准化的事务协议会使集成变得非常困难。有序的事件日志结合幂等的事件处理（参见第11章“幂等性”）是一种更简化的抽象，因此在异构系统中实现时会更可行<sup>[7]</sup>。

基于日志的集成的一大优势是各个组件之间的松耦合，这体现在两个方面：

1. 在系统级别，异步事件流使整个系统在应对各个组件的中断或性能下降时表现更加稳健。如果消费者运行缓慢或失败，那么事件日志可以缓冲消息（请参阅第11章“磁盘空间使用情况”），这样生产者和其他消费者可以继续运行不受影响。





- 有问题的用户恢复正常后也可以继续运行，不会丢失任何数据，这样有效抑制了故障。相比之下，分布式事务的同步交互往往会将本地故障升级为大规模故障（请参见第9章“分布式事务的限制”）。
2. 在人员角度看，分离式数据系统使得不同的团队可以独立的开发、改进和维护不同的软件组件和服务。专业化使得每个团队都可以专注于做好一件事情，且与其他系统维护清晰明确的接口。事件日志提供了一个足够强大的接口，不但能捕获相当强的一致性（通过持久性和事件顺序），同时也普遍适用于几乎任何类型的数据。

## 分离式与集成式系统

即使分离式确实能成为未来的主流，它也不会取代目前形式的数据库，这些数据库仍然会像以前一样拥有巨大需求。维护流处理器中的状态仍然需要数据库，为批处理和流处理器的输出提供合并的查询（请参阅第10章“批处理工作流的输出”和第11章的“流处理”）也同样需要数据库。专门的查询引擎对于特定的工作负载仍然很重要：例如，MPP数据仓库中的查询引擎针对探索性分析查询进行了优化，非常适合处理这种类型的工作负载（请参阅第10章“对比Hadoop与分布式数据库”）。

运行多个不同的基础架构所带来的复杂性可能确是一个问题：每一套软件都需要一个学习曲线，配置问题以及操作习惯等，因此需要尽可能地减少所部署的组件。与这种包含多个组件并依靠应用层代码组合在一起的系统相比，单个集成的软件产品有可能确实在其针对的负载上表现更好，性能更可预测。所以，正如我在序言中所讲的那样，构建你并不需要的扩展规模是在浪费精力，可能会把你限制在一个不灵活的设计中。实际上，这是在做过早优化。

分离的目标不是要与那些针对特定负载的单个数据库来竞争性能。目标是让你可以将多个不同的数据库组合起来，以便在更广泛的工作负载范围内实现比单一软件更好的性能。这里的关注点是广度而非深度。与我们在第10章的“对比Hadoop与分布式数据库”中讨论存储和处理模型的多样性一脉相承。

因此，如果有一种技术可以满足你的所有需求，那么最好使用该产品，而不是试图用基础组件来重新实现它。当没有单一的软件能满足所有需求时，分离和组合的优势才会显现出来。

## 遗漏了什么？

组合数据系统的工具正在变得越来越好，但是我认为还缺少一个主要部分：我们还没





有与UNIX shell相媲美的分离型数据库（即，以简单和声明式构建存储和处理系统的高级语言）。

例如，如果可以简单地声明`mysql | elasticsearch`，类似UNIX管道<sup>[22]</sup>，我一定会非常喜欢这样的实现方式，因为这将是CREATE INDEX的分离式等价物：它将使用MySQL数据库中的所有文档，并将其索引到Elasticsearch集群中。然后，它会不断捕获对数据库所做的更改，并自动将其应用于搜索索引，而无需编写自定义应用程序代码。这种整合应该对几乎所有类型的存储或索引系统都是可行的。

同样，如果能够支持预先计算并方便地更新缓存就更好了。回想一下，实体化视图本质上是一个预计算的缓存，因此可以为复杂的查询以声明式创建实体化视图来缓存结果，包括图的递归查询（参阅第2章“类图数据模型”）以及一些应用处理逻辑。这方面有一些有趣的早期研究，比如差异化数据流<sup>[24,25]</sup>，我希望这些想法能够投入到生产系统中去。

## 围绕数据流设计应用系统

通过应用层代码来组合多个专门的存储与处理系统并实现分离式数据库的方法也被称为“数据库由内向外”方法<sup>[26]</sup>，这是我在2014年的一次会议报告中使用的标题<sup>[27]</sup>。然而，把它称为“新架构”有点夸大了。我把它看作是一个设计模式，一个讨论的起点，名字的目的只是为了方便讨论。

这些想法不是我的，而是一些值得学习的他人的想法的融合。尤其是，它与数据流语言如Oz<sup>[28]</sup>和Juttle<sup>[29]</sup>，功能性反应式编程语言（Functional Reactive Programming, FRP）如Elm<sup>[30,31]</sup>，逻辑编程语言如Bloom<sup>[32]</sup>等存在很多共同点。在此背景下，由Jay Kreps<sup>[7]</sup>提出了**unbundling**这个术语。

甚至电子表格也比大多数主流编程语言更早支持数据流编程功能<sup>[33]</sup>。在电子表格中，可以在一个单元格中放入公式（例如，另一列中的单元格的总和），并且每当公式的任何输入发生更改时，都会自动重新计算公式的结果。这正是我们在数据系统级别所需要的：当数据库中的记录发生更改时，我们希望自动更新该记录的任何索引，并且自动刷新依赖于该记录的任何缓存的视图或聚合。你不必担心技术细节是如何实现这种刷新的，只需要能够确信它可以正常工作即可。

因此，我认为绝大多数的数据系统仍然可以从VisiCalc在1979年已经具备的功能中学习<sup>[34]</sup>。与电子表格的不同之处在于，今天的数据系统需要具备容错性、可扩展性以及







数据持久存储能力。它们还需要能够集成不同团队编写的不同技术，并重用现有的库和服务：期望使用某种特定的语言、框架或工具来开发所有软件是不切实际的。

在本节中，我将对这些想法进行扩展，并围绕着分离式数据库和数据流思路来探索构建新型应用程序的一些方法。

## 应用程序代码作为派生函数

当某个数据集从另一个数据集派生而来时，它一定会经历某种转换函数。例如：

- 二级索引是一种派生的数据集，它具有一个简单的转换函数：对于主表中的每一行或者一个文档，挑选那些索引到的列或者字段值，并且按照值进行排序（假设一个B-tree或者SSTable索引，按主键排序，如第3章所述）。
- 通过各种自然语言处理函数（如语言检测、自动分词、词干或词形识别、拼写检查和同义词识别）创建全文搜索索引，然后构建用于高效查找的数据结构（例如反向索引）。
- 在机器学习系统中，可以考虑通过应用各种特征提取和统计分析功能从训练数据中导出模型。当应用与新输入数据时，模型的输出是基于输入数据和模型（因此间接地从训练数据）派生而来。
- 缓存通常包含那些即将显式在用户界面（UI）的聚合数据。因此填充缓存需要知道在UI中引用哪些字段，UI更改可能也需要相应的调整缓存填充方式和重建方式。

为二级索引构造派生函数的需求非常普遍，作为核心特性，几乎都内置于很多数据库中，通过简单地调用CREATE INDEX即可使用。对于全文索引，常见的语言方面的函数也会内置于数据库中，但更复杂的特性通常还需要针对特定领域进行调整。而对于机器学习，众所周知特征工程与特定应用紧密相关，此外，还需要集成很多用户交互和应用部署方面的很多详细知识<sup>[35]</sup>。

如果创建派生数据集的函数并非创建二级索引那样标准函数时，就需要引入自定义的代码来处理特定应用相关。而这个自定义代码是让很多数据库纠结的地方。虽然关系型数据库通常都支持触发器、存储过程以及用户定义的函数，这些功能使得数据库可以执行应用程序代码，但它们在数据库设计中基本都是后来所添加的功能（请参阅第11章“发送事件流”）。





## 应用程序代码与状态分离

理论上讲，数据库可以像操作系统那样成为任意应用程序代码的部署环境。然而，实际上它们并不太适合，主要是无法满足当今应用程序的很多开发要求，如依赖性和软件包管理，版本控制，滚动升级，可演化性，监控，指标，网络服务调用以及与外部系统的集成。

另一方面，诸如Mesos，YARN，Docker，Kubernetes等用于部署和集群管理的工具是为专门运行应用程序代码而设计的。由于专注于做好一件事，它们能够做得比数据库更好，在众多特性中，都提供了用户定义函数的执行功能。

我认为系统的某部分专注于持久性数据存储，同时有另外一部分专门负责运行应用程序代码是有道理的。这两部会有交互，但是各自仍保持独立运行。

现在大多数Web应用程序都被部署为无状态服务，用户请求都可以被路由到任意的应用程序服务器上，并且服务器发送响应之后就不会保留请求的任何状态信息。这种方式部署起来很方便，服务器可以随意添加或删除，然而有些状态势必需要保存在某个地方：通常是数据库。目前的趋势是将无状态应用程序逻辑与状态管理（数据库）分开：不把应用程序逻辑放在数据库中，也不把持久状态放在应用程序中<sup>[36]</sup>。正如函数程序社区的人喜欢开玩笑说，“我们相信可变状态与Church的工作分离”<sup>[37]注1</sup>。

在这种典型的Web应用模型中，数据库充当一种可以通过网络同步访问的可变共享变量。应用程序可以读取或更新变量，数据库负责持久性，提供一些并发控制和容错功能。

但是，在大多数编程语言中，无法订阅可变变量的更改信息，而只能定期不断地读取它。与电子表格不同，如果变量的值发生了变化，变量的读者将不会收到通知（你可以在你自己的代码中实现这样的通知，例如通过所谓的观察者模式，但是大多数编程语言级别没有内置该模式）。

数据库继承了这种被动方法来处理可变数据：如果想知道数据库的内容是否了变化，唯一的选择就是轮询（即定期地执行查询）。订阅更改只是最近才出现的新功能（请参阅第11章“更改数据流的API支持”）。

---

注1： 解释一个笑话很少能够提高笑话的效果，但我这里解释的目的主要是想让所有人都能懂这个笑话。在这里，数学家Alonzo Church创建了lambda演算，一种早期计算形式后来成为大多数函数式编程语言的基础。Lambda算法没有可变状态（即没有可被覆盖的变量），所以说这是可变状态与Church的工作分离。





## 数据流：状态变化和应用程序代码之间的相互影响

从数据流的角度思考应用意味着重新协调应用代码和状态管理之间的关系。我们不是将数据库简单地视为被应用程序所操纵的被动变量，而是更多地考虑状态、状态变化以及处理代码之间的相互作用和协作关系。应用程序代码在某个地方会触发状态变化，而在另一个地方又要对状态变化做出响应。

我们在第11章的“数据库与数据流”中看到了这样的思路，当时讨论了将数据库更改日志作为可订阅的事件流来处理。消息传递系统（如actors）（请参阅第4章“基于消息传递的数据流”）也具有响应事件的概念。早在20世纪80年代，元组空间模型探索了将分布式计算的过程表示为一些进程负责观察状态的变化，同时一些进程作出响应<sup>[38,39]</sup>。

如前所述，数据库内部也有类似情况，例如当触发器由于数据变化而触发时，或者数据发生变化而二级索引更新时。分离式数据库也采取这样的思路，并将其应用于在主数据库之外创建派生数据集，包括缓存、全文搜索索引、机器学习或分析系统。我们可以使用流处理结合消息传递系统来达到这个目的。

要记住的重要一点是，维护派生数据与异步作业执行不同，后者主要是传统消息系统所针对的目标场景（请参阅第11章“对比日志机制与传统消息系统”相关内容）。

- 当维护派生数据时，状态更改的顺序通常很重要，如果从事件日志中派生出了多个视图，每个视图都需要按照相同的顺序来处理这些事件，以使它们互相保持一致。如第11章“确认与重传”中所述，许多消息代理在重新发生未确认的消息时不支持该特性。而双重写入也被排除在外（请参阅第11章“保持系统同步”）。
- 容错性是派生数据的关键：丢失哪怕单个消息都会导致派生数据集永远无法与数据源同步。消息传递和派生状态更新都必须可靠。例如，默认情况下，许多actor系统默认在内存中维护actor状态和消息，所以如果机器崩溃，这些内存信息都会丢失。

对稳定的消息排序和可容错的消息处理的要求都非常严格，但是它们比分布式事务代价小很多，并且在操作上也更加稳健。现代流式处理可以对大规模环境提供排序和可靠性保证，并允许应用程序代码作为stream operator运行。

这个应用程序代码可以完成那些数据库内置函数所不支持的任意处理逻辑，就像管道链接的UNIX工具一样，多个stream operators可以组装在一起构建大型数据流系统。每个operator以变化的状态为输入，并产生其他状态变化流作为输出。





## 流式处理与服务

当前流行的应用程序开发风格是将功能分解为一组通过同步网络请求（例如REST API）进行通信的服务（请参阅第4章“基于服务的数据流：REST和RPC”）。这种面向服务的结构优于单体应用程序之处在于松耦合所带来的组织伸缩性：不同的团队可以在不同的服务上工作，这减少了团队之间的协调工作（只要服务可以独立部署和更新）。

将stream operator组合成数据流系统与微服务理念有很多相似的特征<sup>[40]</sup>。但是，底层的通信机制差异很大：前者是单向、异步的消息流，而不是同步的请求/响应交互。

除了第4章“基于消息传递的数据流”中所列出的优点（比如更好的容错性）之外，数据流系统还可以获得更好的性能。例如，假设客户正在购买一种商品，这种商品以某一种货币定价，但需要以另一种货币支付。为了执行货币转换，需要知道当前的汇率。这个操作可以通过两种方式来实现<sup>[40,41]</sup>：

1. 对于微服务方法，处理购汇的代码可能会查询汇率服务或数据库，以获取特定货币的当前汇率。
2. 对于数据流方法，处理购汇的代码会预先订阅汇率更新流，并在当地数据库发生更改时记录当前的汇率。当有购汇请求时，只需查询本地数据库即可。

第二种方法把向另一个服务发起的同步网络请求转换为本地数据库的查询（同一个进程或者同一台机器上<sup>注2</sup>）。数据流不仅方法更快，而且当另一个服务失败的话也更加健壮。最快和最可靠的网络请求就是根本没有网络请求！我们现在在购汇请求和汇率更新事件之间有一个事件流join，而不是RPC，请参阅第11章“流和表join操作”。

注意，这种join效果会随时间而变：如果购汇事件在稍后被重新处理，汇率已经发生改变。如果想重演当初的购买行为，则需要从原来购买的时间获得历史汇率。无论是查询服务还是订阅汇率更新流，都需要处理这种时间依赖性（请参阅第10章“join的时间依赖性”）。

订阅变化的流，而不是在需要时去查询状态，使我们更接近类似电子表格那样的计算模型：当某些数据发生更改时，依赖于此的所有派生数据都可以快速更新。还有很多开放的问题，例如关于时间依赖join等，但我相信围绕数据流来构建应用程序是一个非常有前途的方向。

---

注2：在微服务方法中，处理购汇的服务可以通过缓存汇率来避免同步网络请求。但是，为了使缓存保持最新，需要定期轮询更新的汇率，或订阅更改流，后者正是采用了数据流的方法。



## 观察派生状态

概括来讲，上一节所讨论的数据流系统主要侧重如何创建派生数据集（如搜索索引，实体化视图和预测模型）并使其保持最新状态。我们可以把这个过程称为写路径：只要有信息写入系统，它就可能经历批处理和流处理的多个阶段，最终每个派生数据集都会更新以包含新写入的数据。图12-1展示了更新搜索索引的示例。

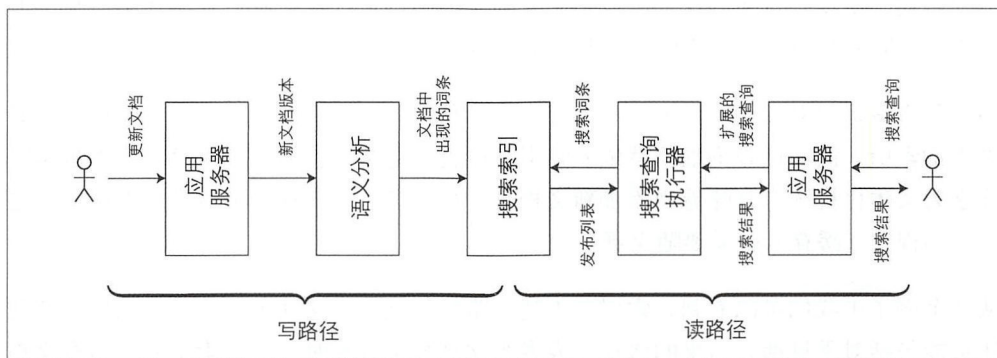


图12-1：搜索索引场景中的写入路径（文档更新）和读取路径（稳定查询）

但是，为什么你首先创建派生数据集呢？很可能是因为想在以后多次查询它即读路径：当有用户访问请求时，从派生数据集中读取数据，做些必要的处理，然后返回处理后的结果以响应用户请求。

总而言之，写路径和读路径涵盖了数据的整个过程，从数据收集到数据使用（可能是由另一个人）。写路径可以看作是预计算的一部分，即一旦数据进入，即刻完成，无论是否有人要求访问它。而过程中的读路径则只有当明确有人要求访问时才会发生。如果你熟悉函数式编程语言，可能会注意到写路径类似于尽早求值，读路径类似于惰性求值。

写路径和读路径在派生数据集上交会，如图12-1所示。某种程度上，它是写入时需完成的工作量和读取时需完成的工作量之间的一种平衡。

## 实体化视图和缓存

全文搜索索引就是一个很好的例子：写路径更新索引，读路径搜索关键字索引。读写都需要做一些工作。写入主要更新文档中出现的所有术语的索引条目，而读则需要搜索输入的每个词，并应用布尔逻辑来查找包含所有词（AND运算符）的文档，或每个词（OR运算符）的任何同义词。



如果没有索引，则搜索查询将不得不扫描所有文档（如grep），如果有大量文档，这样做的代价会很高。没有索引意味着写路径上的工作变少了（没有索引要更新），但是在读路径上则需要更多的工作。

另一方面，可以想象一下为所有可能的查询预先计算搜索结果。此时，读路径上的工作量会减少：不需要布尔逻辑，只需要查找查询结果并返回即可。但是，写路径代价会非常高，主要是可能的查询组合是无限的，因此预先计算所有可能的搜索结果需要几乎无上限的时间和空间代价，这明显实际不可行<sup>注3</sup>。

另一种方案是只对那些最常见的一组查询进行预计算，这样可以快速得到响应而不必去查询索引，而其他查询则仍然依靠索引。这就是通常所称的查询缓存，或者也可以称之为实体化视图，意味着当添加新文档时，如果其包含在那些最常见的查询结果之中，则视图（缓存）就需要随之更新。

从这个例子中我们可以看到，索引并不是写路径与读路径之间链接的唯一边界。其他还可能包括对常见搜索结果的缓存；或者当文档数量有限时，无需索引而对所有文档进行grep方式的全文扫描。从这个角度来看，缓存、索引和实体化视图的角色就比较清楚了，它们主要是调整读、写路径之间的边界。通过预先计算结果，写路径上承担了更多的工作，而读路径则可以简化加速。

读写路径上工作量的转移或平衡实际上正是本书开始时讨论Twitter例子的主题，详见第1章“描述负载”。在该例子中，我们还看到对于那些明星账户，读写路径可能与普通用户有很大差异。再次回顾，我们已经完成了本书近几百页的讨论。

## 有状态，可离线客户端

我发现读写路径之间的这种边界非常有趣，可以继续讨论调整边界，以及实际场景中这种调整究竟意味着什么。接下来，让我们看看不同场景下的具体情况。

由于Web应用程序在过去二十年的普及流行，很多人由此对应用开发有了一些先入为主的假设。特别是，C/S架构即客户端/服务器端模型如此普遍（客户端大部分是无状态，服务器端拥有对数据的全部控制权），以至于大家几乎都忽略了其他架构的存在。但是，技术永远在不断发展，我认为经常质疑现状是非常重要的。

传统上，Web浏览器是无状态的客户端，只有当互联网连接时才能触发后面的处理

---

注3： 绝不是开玩笑，假设有一个有限的语料库，明确清晰的非空搜索结果的查询可能是有限的。然而，语料库中的词条数量确是指数级的，这将是一个噩梦。



（离线时，唯一能做的事情就是在之前以打开的网页上滚动浏览）。然而，最近的“单页”JavaScript Web应用程序已经支持很多有状态的功能，包括基于客户端的可交互用户界面和数据持久保存在Web浏览器。移动应用程序类似地可以在设备上保存很多状态，并且很多交互不需要与服务器进行往返通信。

这些新特性重新引发了对离线优先应用程序的兴趣，它们尽可能地在同一设备上使用本地数据库，而不需要互联网连接，并在网络连接可用时在后台与远程服务器同步<sup>[42]</sup>。由于移动设备的互联网连接通常比较慢，而且可靠性低，因此如果用户界面不需要等待同步网络请求，并且大多数应用程序可以离线工作，对用户来说是一个巨大优势（请参阅第5章“离线客户端操作”）。

当我们摆脱无状态客户端与核心数据库这样的假设，而转向在终端用户设备上维护状态时，就打开了一个全新机遇。特别是，我们可以将设备上的状态视为服务器上的状态缓存。屏幕上的呈现是一种客户端对象模型的实体化视图；而客户端的对象模型则是远程数据中心在本地的状态副本<sup>[27]</sup>。

## 状态更改推送至客户端

对于典型的Web网页，在Web浏览器加载页面之后，如果服务器端数据发生变化，则只能重新页面才能观察到新数据。浏览器只在某一个时间点读取数据，假设它是静态的，它不会订阅服务器端的更新。因此，除非显式地轮询更新，设备上的状态是旧的缓存（像RSS这样基于HTTP的订阅协议实际上只是一种轮询形式）。

之后更新的协议已经超越了基本的HTTP的请求/响应模式：服务器端发送事件（EventSource API）和WebSocket提供了一个通信通道，通过该通道，Web浏览器可以与服务器保持开放的TCP连接，只要处于连接状态，服务器可以主动推送消息至浏览器。这为服务器提供了一种新的更新机制，即主动通知并更新终端用户客户端本地存储状态，从而缩小与服务器状态的滞后程度。

就我们的写路径和读路径模型而言，主动推送状态至客户端设备意味着将写路径一直延伸到终端用户。当客户端首次初始化时，仍然需要使用读取路径来获得初始状态，但此后可能依赖于服务器发送的状态更改流。所以，流处理和消息传递方面的这些想法并不局限于仅在数据中心运行，我们可以进一步并将这些想法扩展到终端用户设备<sup>[43]</sup>。

这些设备有时会离线，在此期间无法接收到服务器状态更改的任何通知。但是我们已经了解解决方案，在第11章“消费者偏移量”中，我们讨论了基于日志的消息代理的消费者在失败或断开连接后可以重新连接，并确保它不会错过任何网络断开之后的新消

息。同样的技术适用于个人用户，即每个设备都抽象为小型事件流的一个的订阅者。

## 端到端的事件流

最近的一些针对有状态客户端和用户接口的开发工具，如Elm语言<sup>[30]</sup>和Facebook的React，Flux和Redux工具链，已经支持从服务器端订阅代表用户输入的事件流来管理内部的客户端状态，结构与事件溯源类似（请参阅第11章的“事件溯源”）。

自然地，可以将这种编程模型扩展为允许服务器将状态改变事件推送至客户端事件流水线中。因此，状态变化可以通过端到端的写路径流动：某个设备上交互行为触发了状态变化，通过事件日志、派生数据系统和流式处理等，一直到另一台设备上用户观察到状态。这些状态变化传播的延迟可以做到很低水平，例如端到端只需一秒。

某些应用程序（例如即时消息系统和在线游戏）已经采用了这种“实时”结构（主要是从低延迟的交互意义上说，而不是第8章“响应时间保证”）。那么为什么并非所有应用都以这种方式构建呢？

主要挑战在于，目前的数据库、函数库、开发框架和交互协议等都有对无状态客户端和请求/响应交互根深蒂固的假设。许多数据存储支持的读写操作其实都是一个请求对应一个响应，但很少支持订阅更改，即一段时间内会主动返回一系列的响应（请参阅第11章“更改事件流的API支持”）。

为了将写路径扩展到最终用户，我们需要从根本上重新思考构建这些系统的方式：从请求/响应交互转向发布/订阅数据流<sup>[27]</sup>。我认为更具响应性的用户界面和更好的离线支持的优势是绝对值得尝试的。如果你正在设计数据系统，我希望你能够记住订阅更改这一方式，而不仅仅是查询当前的状态。

## 读也是事件

我们讨论了当流式处理将派生数据写入存储（数据库，缓存或索引）时，以及当用户请求查询存储时，存储充当了上述写路径与读路径之间的一种可调边界。当查询时，存储可以支持对数据进行随机快速访问（例如基于索引等），或者扫描整个事件日志（没有索引时）。

很多情况下，数据存储与流处理系统是分开的。但流式处理自身还需维护一些状态以方便执行聚合和join（参阅第11章“流式join”）。这些状态通常隐藏在流处理系统内部，但是也有一些框架支持外部客户端查询<sup>[45]</sup>，这样流处理系统本身也成为一种简单的数据库。



我想进一步来谈这个想法。正如前面所讨论的，对存储的写入是通过事件日志进行的，而读取则是即时的网络请求方式，查询直接路由到那些存储数据节点。这样的设计也很合理，但它并不是唯一可能的设计方案。也可以将读请求表示为事件流，发送至流处理系统，流处理系统则将读结果发送至输出流来响应读取事件<sup>[46]</sup>。

当写入和读取都被表示为事件，并且被路由到相同的stream operator统一处理时，我们实际上是在查询流和数据库之间执行stream-table join操作。读事件需要发送到保存数据的数据库分区节点上（参阅第6章“请求路由”），这一点与批处理和流处理系统在join操作时需要在同一个主键上合并数据一样（参阅第10章“reducer端的join与分组”）。

服务请求与执行join操作之间的对应关系非常重要<sup>[47]</sup>。一次性读取请求只是通过join运算来传递请求，之后就再也不用；然后订阅请求则是对另一侧所有过去和未来事件的一种持续性的join操作。

以日志方式记录读事件可能还可以帮助跟踪系统级别的事件因果关系和数据源：它可以重建用户在做出某个决定之前看到的内容。例如，一个购物网站，向客户显示的预计出货日期和库存状态可能会影响他们选择购买某件物品<sup>[4]</sup>。为了分析此间的关联，需要记录用户所查询到的相关快递和库存状态。

将读取事件信息写入持久化存储，一方面可以更好地追踪因果关系（参阅本章前面“排序事件以捕获因果关系”），同时会导致额外的存储和I/O开销。对此进行优化以减少额外开销目前仍是一个开放的研究话题<sup>[2]</sup>。但是，如果你已经出于运维目的而将读信息记录到日志，作为请求处理的副作用，附带支持将日志本身也作为请求的分析来源其实并没有改变太多。

## 多分区数据处理

对于仅涉及单个分区的查询，通过流来发送查询并收集响应事件流可能显得有些大材小用。然而，这种方法却开启了一种分布式执行复杂查询的可能性，这需要合并来自多个分区的数据，并很好地借助底层流处理系统所提供的消息路由、分区和join功能。

Storm的分布式RPC功能支持这种使用模式（请参阅第11章“消息传递和RPC”）。例如，用来计算Twitter上有多少人观看了某个URL地址，即推送该条URL信息的所有用户的所有关注者的并集<sup>[48]</sup>。由于Twitter用户采用了分区机制，所以上述计算也需要合并多个分区上的查询结果。





这种模式的另一个例子是反欺诈：为了评估特定购买事件是否具有欺诈风险，可以检查用户的IP地址、电子邮件地址、账单地址、送货地址等等的信誉分数。这些信誉数据库每一个都是独立分区的，因此为特定购买事件收集最终分数需要一系列对不同分区数据的join<sup>[49]</sup>。

MPP数据库的内部查询执行计划也具有相似的特征（请参阅第10章的“对比Hadoop与分布式数据库”）。当需要执行这种跨分区的join操作时，数据库内置提供此功能可能比使用流处理系统实现起来更简单。然而，将查询也作为一种事件流则提供了一个新的选择，来实现大规模应用系统并摆脱现有传统方案的限制。

## 端到端的正确性

对于只读取数据的无状态服务，如果发生错误，这并不是什么大问题：可以修复错误并重新启动服务，然后一切恢复正常。但是像数据库这样的有状态系统就不是那么简单了：它们被设计成永远记住事物（有的多有的少），所以一旦发生问题，影响可能会是永远性的，这就要求必须仔细对待<sup>[50]</sup>。

我们总是希望构建可靠并且正确无误的应用程序（即使在面对各种故障的情况下，其语义依然保持清晰定义和可理解）。在最近四十年里，原子性，隔离性和持久性（第7章）的事务特性一直是保证正确性的首选方式。但是，这些基础特性实际上却比他们看起来要弱，例如我们已经见识了弱隔离级别的很多混淆不清的现状（请参阅第6章“弱隔离级别”）。

在某些领域，事务处理被完全抛弃，被性能更好、扩展性更强的模型所取代，但是这些模型具有更复杂的语义（例如第5章“无主节点复制”）。一致性虽然经常被关注，但始终缺乏明确的定义（参见第7章和第9章的“一致性”）。业界还有声音主张，为了更好地可用性，我们应该“拥抱弱一致性”，然而实践中类似这样说法缺乏清晰的建议。

对于如此重要的话题，我们的理解和手头的工程方法确显得非常薄弱。例如，在特定事务隔离级别或复制模式下，想要确定某个特定应用程序是否安全就很有挑战性<sup>[51,52]</sup>。简单的解决方案似乎能够降低并发且没有错误的情况下正常工作，但是在要求更高的情况下会出现许多细微的错误。

例如，Kyle Kingsbury的Jepsen实验<sup>[53]</sup>已经表明，当存在网络问题和节点崩溃时，某些产品所声称的安全保证与实际行为之间存在明显差异。即使像数据库这样的基础架构产品没有问题，上层应用程序代码仍然需要正确地使用它们所提供的功能（包括弱





隔离级别、quorum配置以及其他），如果配置难以理解，就很容易出错。

如果你的应用程序能够容忍偶尔发生一些不可预知的崩溃或丢失数据，那么情况就简单得多，或许你可以叉着手很轻松的置身事外。否则的话，只能需要更强的正确性保证，即可串行化和原子提交，但它们都是有代价的：通常只能用于单数据中心（因此排除了地域分布的部署架构），有限的扩展规模和容错性。

传统的事务方法并不会消失，但我也相信，它也未必是保证应用程序正确和灵活处理错误的终极手段。本节我将提出一些基于数据流架构来思考、探索正确性的若干方法。

## 数据库的端到端争论

仅仅因为应用程序使用了具有较强安全属性的数据系统（例如可串行化的事务），并不能意味着应用程序一定保证没有数据丢失或损坏。例如，如果应用程序的bug导致它写入不正确的数据，或者从数据库中删除数据，那么可串行化的事务也无济于事。

这个例子看似简单，但这些情况却值得重视，包括应用程序存在bug和人为操作存在失误。在第11章“状态，流和不可变性”中我使用了这个例子来阐述不可变数据和追加式更新，这样如果把错误的代码修复或移除掉，可以更容易从错误中恢复。

尽管不可变性很有帮助，但它本身并不是一种万能的方法。让我们看看可能发生数据损坏的一个更微妙的例子。

### Exactly-once执行操作

在第11章“流处理的容错”中，我们提到了称为“只执行一次”（*exactly-once*，或者*effectively-once*）语义。如果在处理消息过程中出现意外，可以选择放弃（丢弃消息，即导致数据丢失）或者再次尝试。如果采用重试策略，而实际上第一次执行时已经事实成功，只是用户没有得到成功的确认，那么最终这个消息其实是处理了两次。

处理两次是某种形式的数据损坏，例如相同的服务向客户收费两次（超额计费）或计数器增加两次（夸大了一些度量）都是不可取的。在这种情况下，*exactly-once*意味着合理安排计算，使得最终效果与没有发生错误的结果一样，即使操作实际上由于某种故障而被重试。我们以前讨论过实现这个目标的几种方法。

最有效的方法之一是使operator满足幂等性（参阅第11章“幂等性”），即无论执行一次还是多次，确保具有相同的效果。但是，如果操作本身并非幂等，而要改造使其具有幂等性则需要付出一定努力，例如可能需要维护一些额外的元数据（例如已成功





完成的操作ID集合），并确保在节点失效、切换过程中采取必要的fencing措施（参阅第8章“主节点与锁”）。

## 重复消除

除了流处理之外，还有许多其他模式也需要消除重复。例如，TCP使用序列号来检测网络上是否有数据包丢失或重复，并最终确保数据包以正确的顺序接受。丢失的数据包都会被重新发送，并且在将数据交给应用程序之前，TCP堆栈将负责删除重复的数据包。

但是，这种重复消除只能用于单个TCP连接上下文中。假设客户端与数据库之间采用TCP连接，并且正在执行示例12-1中的事务。在许多数据库中，事务与某个客户端连接绑定（如果客户端发送了多个查询，且运行在同一个TCP连接上，则数据库就知道它们属于同一个事务）。如果客户端在发送COMMIT之后，收到响应之前遭遇网络中断和连接超时等，它就不知道事务是否已被提交或中止（见图8-1）。

示例12-1：从一个账户到另一个账户的非幂等性转账

```
BEGIN TRANSACTION;  
UPDATE accounts SET balance = balance + 11.00 WHERE account_id = 1234;  
UPDATE accounts SET balance = balance - 11.00 WHERE account_id = 4321;  
COMMIT;
```

客户端可以重新连接到数据库并重试事务，但现在它不在TCP可重复消除的范围之内。由于示例12-1中的交易不是幂等的，可能会发生转了22美元，而不是实际希望的11美元。因此，尽管示例12-1是一个标准的交易原子性例子，但它实际上是有些问题，真实的银行不是像这样工作<sup>[3]</sup>。

两阶段提交（请参阅第9章“原子提交与两阶段提交”）协议会破坏TCP连接和事务之间的1:1映射关系，因为在网络故障之后，事务协调者必须重新连接数据库，来通知事务是决定提交还是有问题因而中止。这是否能够确保交易只能执行一次吗？不幸的是，答案是否定的。

即使可以消除数据库客户端与服务器之间的重复事务，我们还得关注最终用户的设备和应用程序服务器之间的网络情况。比如下面这种情况，如果最终用户客户端是Web浏览器，则可能使用HTTP POST请求向服务器提交指令。也许用户在一个缓慢的移动蜂窝线路上，虽然成功地发送了POST，但是当服务器响应时网络信号太弱以至于无法成功接收。

在这种情况下，用户端很可能会显示一条错误消息，并可能不得不手动重试。Web浏览器警告说“确定要再次提交表单吗？”，用户单击是的，他希望看到操作成功的提







示（Post/Redirect/Get模式<sup>[54]</sup>可以避免正常操作中的类似警告消息，但是如果POST请求超时，这种模式也将无济于事）。从Web服务器的角度来看，重试是一次独立的请求，而对于数据库来说重试请求变成一个另外全新的事务。通常的去重机制对此无能为力。

## 操作标识符

为了实现跨多次网络跳转请求而操作仍然具有幂等性，仅仅依靠数据库提供的事务机制是不够的，需要考虑请求的端到端过程。

例如，可以为操作生成一个唯一的标识符（如UUID），并将其作为一个隐藏的表单字段包含在客户机应用程序中；或者对所有相关表单字段计算一个哈希值依次来代表操作ID<sup>[3]</sup>。如果Web浏览器两次提交POST请求，则两个请求具有相同的操作ID。然后，可以将该操作ID一直传递到数据库，检查并确保对一个给定的ID只执行一个操作，参见示例12-2。

示例12-2：采用唯一的ID来消除重复的请求

```
ALTER TABLE requests ADD UNIQUE (request_id);

BEGIN TRANSACTION;

INSERT INTO requests
(request_id, from_account, to_account, amount)
VALUES('0286FDB8-D7E1-423F-B40B-792B3608036C', 4321, 1234, 11.00);

UPDATE accounts SET balance = balance + 11.00 WHERE account_id = 1234;
UPDATE accounts SET balance = balance - 11.00 WHERE account_id = 4321;

COMMIT;
```

示例12-2依赖于`request_id`列上的唯一性约束。如果一个事务试图插入一个已存在的ID，那么INSERT失败，事务被中止而退出，这样无法执行两次。即使在弱隔离级别下，关系型数据库通常也可以正确地确保唯一性约束，而“应用程序采用先检查-再插入”的方式可能会在不可串行化隔离下失败，参见第7章“写倾斜与幻读”的例子。

除了消除重复请求之外，示例12-2中`requests`表作为一种事件日志，还可以表征事件源（请参阅第11章“事件溯源”）。账户余额的更新操作不一定要和插入处于同一个事务内，它们可以是冗余的，从下游消费者的请求事件中派生而来，而基于唯一请求ID可以确保事件被处理一次。

## 端到端的争论

消除重复事务这种场景其实只是一种更为普遍的“端到端论点”原则的其中一例，该





原则最早由Saltzer、Reed和Clark与1984年所阐述<sup>[55]</sup>：

只有具备应用程序充分的知识，并且站在通信系统端点的角度的情况下，才能完全正确地实现所关注的功能。因此，以通信系统本身的特征来提供这种被质疑的功能是不可能的（有时，由通信系统提供的不完整版本则可能有助于提高性能）。

对于我们的例子，所关注的功能是重复消除。我们看到TCP协议可以在连接层消除重复的数据包，而一些流处理系统在消息处理级别提供了所谓的exactly-once语义，但是这还不足以防止用户在第一次超时后提交重复的请求。TCP，数据库事务和流处理器本身并不能完全排除这些重复。解决这个问题需要一个端到端的解决方案：从终端用户客户端一直传递到数据库的唯一事务标识符。

端到端的参数也适用于检查数据的完整性：以太网，TCP和TLS中内置的校验码可以检测网络中数据包的损坏情况；但是由于网络收发两端的软件bug或存储在磁盘上的数据损坏则无法检测到。如果想要捕获所有可能的数据源损坏，则还需要端到端的校验和。

类似的讨论也适用于加密<sup>[55]</sup>：家庭WiFi网络上的密码可以防止人们窃听WiFi流量，但无法防范联网过程中其他地方的攻击；而客户端和服务端之间的TLS / SSL加密则可以防止整个网络攻击，但不会保护服务器本身的安全。只有端到端的加密和认证可以防止所有这些事情。

尽管这些底层功能（TCP复制消除，以太网校验和，WiFi加密）无法单独提供所需的端到端功能，但它们仍然有用，因为它们可以降低更高层出现问题的可能性。例如，如果TCP层不负责保证正确的数据包顺序，那么HTTP请求就可能出现混乱。我们只需要记住，底层的可靠性功能本身不足以确保端到端的正确性。

## 在数据系统中采用端到端的思路

这里我要强调一下我原来的观点：即使应用程序所使用的数据系统提供了比较强的安全属性（如可串行化事务），也并不意味着应用程序就一定没有数据丢失或损坏，应用程序本身也需要采取端到端的措施，例如重复消除。

这真是一个耻辱，时至今日，容错机制还是很难保证最后正确的结果。底层的可靠性机制（如TCP中的可靠性机制）没有问题，由此已经降低了高层故障的发生概率。因此，想办法把其余高层的容错机制封装成一个抽象层听起来是个不错的主意，这样就不必由应用程序代码担心了，然而我担心目前还没有找到这样正确的抽象层。





长期以来，事务被认为是一个非常好的抽象，我相信它的确非常有用。正如在第7章所讨论的那样，它把可能出现的诸多问题（并发写入，违反约束，崩溃，网络中断，磁盘故障等）归结为只有两种可能的结果：提交或中止。这对编程模型来讲是一个巨大的简化，但我担心这依然不够。

事务处理的代价很高，特别是在涉及异构存储技术时（请参阅第9章“实践中的分布式事务”）。当我们因为代价昂贵而拒绝采用分布式事务，最终不得不在应用代码中重新实现容错机制。正如本书中大量的例子所揭示的，关于并发性和部分错误的推理非常困难而且不够直观，所以我谨慎怀疑大多数应用程序级别的容错机制可能不能正常工作。最后结果是丢失或损坏数据。

基于以上种种考虑，我认为探索更好的容错抽象非常必要：它能够更容易地提供特定应用的端到端正确保证，并且在大规模的分布式环境中依然具有良好的性能和良好的操作性。

## 强约束

让我们探讨一下分离式数据库背景下的正确性（本章前面的“分离式数据库”）。我们看到，使用从客户端一直传递到数据库的请求ID，可以实现端到端的重复消除。那么其他类型的限制怎么样呢？

我们特别关注唯一性约束，即示例12-2所依赖的约束。第9章“约束与唯一性保证”还有更多需要保证唯一性的应用示例，包括用户名或电子邮件地址必须唯一标识用户，文件存储服务不能有多个相同名称的文件，两个人不能预定同一个航班或剧院相同的座位等。

其他类型的约束条件非常类似，例如，确保账户余额永远不会变为负数，不能出售多于仓库库存的商品，或者会议室不能重复预订等。保证唯一性的技术通常也可以用于这类约束。

## 唯一性约束需要达成共识

在第9章中我们介绍了在分布式环境中，保证唯一性约束需要达成共识：如果有多个具有相同值的并发请求，系统需要决定接受哪一个操作，并由于违法约束因此拒绝其他的冲突操作。

达成这一共识的最常见的方式是将单一节点作为主节点，并负责作出所有的决定。只要不介意通过单个节点发送所有请求（即使客户端来自地球另一边），并且只要该节





点不出问题，就能够达成共识。如果需要容忍主节点出错，那么又重新回到了共识问题上（请参阅第9章“主从复制与共识”）。

采用分区方案可以提高唯一性检查的扩展性，即基于要求唯一性的字段进行分区。例如，如果需要通过请求标识来确保唯一性，参见示例12-2，可以把相同请求标识的所有请求都路由到同一分区（请参阅第6章）。如果需要用户名唯一，可以通过用户名的哈希值进行分区。

但是，它无法支持异步的多主节点复制，因为可能会发生不同的主节点同时接受冲突写入，无法保证值的唯一（请参阅第9章“实现线性化系统”）。如果想达到立即拒绝任何违反约束的写请求，则同步式的协调不可避免<sup>[56]</sup>。

### 基于日志的消息传递唯一性

日志机制可以确保所有消费者以相同的顺序查看消息，这种保证在形式上被称为全序关系广播，它等价于共识问题（详见第9章“全序关系广播”）。在基于日志的消息传递的分离式数据库系统中，我们可以采用非常类似的方法来保证唯一性约束。

流处理系统在单线程上严格按照顺序来消费处理日志分区中的所有消息（请参阅第11章“对比日志机制与传统消息系统”）。因此，如果根据需保证唯一性的字段进行日志分区，则流处理系统可以清晰、明确地确定多个冲突操作哪一个先到达。我们来看一个典型的例子：多个用户尝试声明相同用户名<sup>[57]</sup>：

1. 设置用户名的每个请求都编码为一条消息，并追加到根据用户名哈希值来确定的分区。
2. 一个流处理操作按顺序读取日志中的请求，基于本地数据库来跟踪记录已经使用了哪些用户名。对于用户名可用的每一个请求，它记录下所取得的用户名，并向输出流发出一条成功消息。对于那些用户名已被占用的请求，则向输出流发出拒绝消息。
3. 请求用户名的客户端观察输出流，并等待请求是否成功（或拒绝）。

这个算法与第9章“采用全序关系广播实现线性化存储”基本相同。通过增加分区数量，可以轻松扩展来应对较大的请求吞吐量，这样每个分区都可以独立处理。

这个方法不仅适用于唯一性约束，而且适用于许多其他类型的约束。其基本原理是，任何可能冲突的写入都被路由到特定的分区并按顺序处理。正如第5章“什么是冲突”和第7章“写倾斜与幻读”所述，冲突的定义可能取决于应用程序，但流处理系



统可以使用任意逻辑来验证请求。这个想法与Bayou在20世纪90年代所开创的方法类似<sup>[58]</sup>。

## 多分区请求处理

当涉及多个分区时，确保操作原子执行且同时满足各种约束是件很有趣的事情。在示例12-2中，可能有三个分区：一个包含请求ID，一个包含收款人账户，另一个包含付款人账户。没有理由一定要把他们三个放在一个分区里，因为本来它们就是相互独立的。

在传统的数据库方法中，执行这个事务需要横跨所有三个分区进行原子提交，这实质上是对三个分区上的所有事物执行全序排列。由于要求跨分区协调，不同的分区不能再独立处理，吞吐量因此会收到显著影响。

但是，事实证明，使用分区日志是可以实现同等的正确性，并且不需要原子提交：

1. 账户A向账户B转账的请求由客户端赋予一个唯一的请求ID，基于该请求ID追加到对应的日志分区。
2. 流处理系统读取请求日志。对于每个请求消息，它发出两条输出消息：到付款人账户A的付款指令（按照A进行分区），以及到收款人账户B（按照B进行分区）的信用指令。原始的请求ID都包含在这两条消息中。
3. 后续操作接收上述指令，并通过请求ID进行重复数据消除，将最后的更改应用于账户余额。

步骤1和步骤2是必要的，因为如果客户端直接发送信用和付款指令，则需要在这两个分区之间进行原子提交以确保两者都成功（或者都不成功）。为了避免处理分布式事务，我们首先持久地将请求记录为单条消息，然后从第一条消息中派生出信用和付款两条指令。单一对象写入在几乎所有的数据系统中都是原子的（请参阅第7章“单一对象写入”），因此请求既可以通过日志消息，或者其他非日志方式，最终避免多分区的原子提交。

如果步骤2的流处理操作发生崩溃，则从上一个快照检查点来恢复处理。这样做，就不会跳过任何请求消息，但它可能会多次处理请求，并产生重复的信用和付款指令。但是，由于它是确定性的，它只会再次产生相同的指令，步骤3中的处理可以基于端到端的请求ID来轻松实现重复数据消除。

如果还想确保付款人账户不会发生透支，可以添加另一个流处理操作符（按照付款人



账号进行分区)维护账户余额并验证转账金额。在步骤1中,只有有效的事务才能放入到请求日志中。

通过将多分区事务划分为两个不同分区的处理阶段,并使用端到端的请求ID,我们实现了同样的正确性(每个转账请求只对付款人和收款人账户应用一次),即使可能会发生故障,而且避免使用原子提交协议。采用多个不同分区的处理阶段这一想法类似于我们在本章“多分区数据处理”一节所讨论的内容(另请参阅第11章“并发控制”)。

## 时效性与完整性

事务的一个非常方便的属性是它们通常是可线性化的(请参阅第9章“线性化”):即,一个写入者等待直到事务提交,提交一旦完成,写入值立即对所有读者可见。

而对于跨多个处理阶段的流操作处理,则情况并非如此:日志的消费模式是基于异步设计的,因此发送者不会等待消息处理。但是,客户端可能会等待消息出现在输出流中。例如,需要检查是否满足唯一性约束时,即本章前面“基于日志消息传递的唯一性”所讨论的。

在上述例子中,唯一性检查的正确性不取决于消息的发送者是否等待结果。等待只是为了同步的通知发送者唯一性检查是否成功,而该通知可以与处理消息剥离开。

概括来讲,我认为一致性这个术语将两个值得分开考虑的不同的需求:时效性和完整性合二为一了:

### 时效性

时效性意味着确保用户观察到系统的最新状态。之前我们看到,如果用户从数据的旧副本中读取数据,他们可能会观察到不一致的状态(请参阅第5章的“复制滞后”)。但是,这种不一致是暂时的,最终通过等待和再次尝试来解决。

CAP理论(请参阅第9章“线性化的代价”)提供基于线性化的一致性,这是实现时效性的强有力的方法。弱时效性,比如`read-after-write`一致性(请参阅第5章“读自己的写”)也很有用。

### 完整性

完整性意味着避免数据损坏,即没有数据丢失,也没有互相矛盾或错误的数据。尤其是,如果将某些派生数据集作为基础数据的视图来进行维护(请参阅第11章“从事件日志导出当前状态”),派生必须做到正确。例如,数据库索引必须正确反映数据库的内容,否则部分缺失的索引实际用处大打折扣。



如果完整性受到破坏，这种不一致将是永久性的：在大多数情况下，等待并再次尝试不是来修复数据库损坏的。相反，数据库需要专门的检查和修理。在ACID事务的上下文中（参见第7章“ACID的含义”），一致性通常被理解为某种特定于应用的完整性概念。原子性和持久性则是保持完整性的重要手段。

简而言之，违反时效性导致“最终一致性”，而违反完整性则是“永久性不一致”。

我敢说对于大多数应用，完整性比时效性重要得多。违反时效性可能令人讨厌和混淆，但是对完整性的破坏则是灾难性的。

例如，在信用卡对账单上，如果你在过去24小时内完成的交易尚未出现，这并不奇怪，这些系统有一定的滞后性是正常的。我们知道银行都是异步的协调和解决结算交易的，时效性在这里并非头等重要<sup>[3]</sup>。然而，如果报表余额不等于交易金额加上之前的报表余额（金额上出现的错误），或者有些交易向你扣款完成但却没有支付给商家，后果将非常严重（资金凭空消失）。这样的问题属于严重破坏系统的完整性。

## 数据流系统的正确性

ACID事务通常既提供时效性（例如线性化）保证，又提供完整性（例如原子提交）保证。因此，如果从ACID事务角度来看待应用程序的正确性，那么时效性和完整性之间的区别其实无关紧要。

另一方面，本章所讨论的基于事件的数据流系统一个有趣特性是它将时效性和完整性分开了。在异步处理事件流时，除非在返回之前明确地创建了等待消息到达的消费者，否则不能保证及时性。但是，完整性仍是流处理系统的核心。

只执行一次（*exactly-once*或*effectively-once*）（请参阅第11章“流处理的容错”）是一种保证完整性的机制。如果事件丢失或者事件发生两次，数据系统的完整性可能被破坏。因此，容错的消息传递和重复消除（例如，幂等操作）对于面对故障时保持数据系统的完整性非常重要。

正如我们在上一节看到的那样，可靠的流处理系统可以在不需要分布式事务和原子提交协议的情况下保持完整性，这意味着它们可以实现等价的正确性，同时具有更好的性能和操作鲁棒性。我们主要是通过以下机制实现了这一完整性：

- 将写入操作的内容表示为单条消息，可以轻松地采用原子方式编写，这种方法非常适合事件源（请参阅第11章“事件溯源”）。
- 使用确定性派生函数从该条消息派生所有其他状态的更新操作（与存储过程类

似）（请参见第7章“实际串行执行”和本章前面的“应用程序代码作为派生函数”）。

- 通过所有这些级别的处理来传递客户端生成的请求ID，实现端到端重复消除和幂等性。
- 消息不可变，并支持多次重新处理派生数据，从而使错误恢复变得更容易（请参阅第11章“不可变事件的优势”）。

在我看来，上述多种机制的组合是构建未来容错应用的一个非常有前途的方向。

## 宽松的约束

如前所述，保证唯一性约束需要达成共识，通常通过单个节点汇集特定分区的所有事件来实现。如果想要达到传统的唯一性约束，上述处理限制对于流处理系统是不可避免的。

然而，另一个需要注意的是，许多实际应用程序采用了较弱的唯一性因而可以摆脱该限制：

- 如果万一两个人同时注册了相同的用户名或预订了同一个座位，则可以向其中一个发送道歉消息，并要求他们选择另一个。这种纠正错误的措施被称为补偿性事务<sup>[59,60]</sup>。
- 如果客户订购的商品超出当前库存，则可以追加补充库存，但需要为延误发货向客户道歉，并为他们提供折扣。实际上，如果叉车运走的仓库中的一些物品超过了限额，实际留下的货品也会比想象的还要少一些，这 and 用户超售最终效果是一样的<sup>[61]</sup>。因此，道歉的工作流程也需要成为业务流程的一部分，因此可能就没有必要对库存余量进行线性化约束了。
- 同样地，许多航空公司会超售机票，因为它们预计有些乘客会错过航班，许多酒店也会超订房间，因为他们也预计有些客人会取消预定。在这种情况下，出于商业原因严格违反了“每个座位一人”的约束，并且在需求超过供应时要处理补偿过程（退款，升级，在邻近酒店提供免费房间等）。即使没有超额预订，为了应对由于恶劣天气而取消的航班或职员罢工等，也需要道歉和赔偿程序。出现这些问题而进行恢复属于业务的正常组成部分<sup>[3]</sup>。
- 如果万一有人提款额比他们账户中的钱还多，银行可以向他们收取透支费用，并要求他们偿还欠款。通过限制每天的提款总额，银行的风险可以控制在限定范围的。

因此，在许多商业环境中，实际上可以接受的是暂时违反约束，稍后通过道歉流程来修复。道歉（金钱或声誉）的成本各不相同，但通常很低：你不能取消已发送的电子邮件，但可以发送后续电子邮件并进行更正。如果不小心向信用卡收取了两次费用，可以退还其中一笔费用，而你的成本只是处理成本，或者客户投诉。如果一个账户透支了，客户也不还款，原则上银行可以通过收债员来追回款项，但是通过自动提款机超额取款，则没法直接把这笔钱收回来。

道歉的成本是否可以接受是一个商业决策。如果可以接受的话，在写入数据之前检查所有约束的传统模型就是不必要的限制，并且不需要线性化的约束。继续写入操作，并在既成事实之后检查约束，可能就成为合理的选择。你仍然可以在发生恢复代价昂贵的事情之前进行验证，但并不意味着在写入数据之前必须进行验证。

而这些应用程序都非常需要完整性：你不会希望失去预留，或者由于错误匹配的信用和借记信息而使资金消失。但是他们并不需要时效性来执行这个约束：如果卖出的物品多于库存量，那么可以通过道歉来事后补救问题。这些做法与第5章“处理写入冲突”中讨论的冲突解决方法类似。

## 无需协调的数据系统

现在我们有了两个有趣的观察：

1. 数据流系统可以保证派生数据的完整性，无需原子提交，线性化或跨分区的同步协调。
2. 虽然严格的唯一性约束要求时效性和协调性，但是只要整体上保证完整性，即使发生暂时约束破坏，可以事后进行修复，因此许多应用实际上采用宽松式的约束并没有问题。

总之，这些观察意味着数据流系统可以为许多应用程序提供数据管理服务，而不需要协调，同时仍然提供强大的完整性保证。这种避免协调的数据系统具有很大的吸引力：与需要执行同步协调的系统相比，可以实现更好的性能和容错能力<sup>[56]</sup>。

例如，它可以用于跨多个数据中心的多主节点复制系统，在区域之间异步复制。由于不需要跨区域同步协调，任何一个数据中心都可以独立运行。虽然时效性保证很弱（如果不引入协调者就不可能实现线性化），但是它仍然具有很强的完整性保证。

在这种情况下，可串行化事务作为维护派生状态的一部分还是有用的，但是只能在一个限定的范围内工作很好<sup>[8]</sup>。异构分布式事务（如XA事务）并不是必需的（请参阅第9章“实践中的分布式事务”）。同步协调可以仅在需要的地方才引入（例如，





在不可能恢复的操作之前执行严格的限制），但是如果只有一小部分应用程序需要它，就没有必要所有事务都进行协调<sup>[43]</sup>。

另一种理解协调和约束的方法是：它们减少了由于不一致而引发的道歉数量，但是也可能降低系统的性能与可用性，并由此可能增加由于业务中断而引发的道歉数量。你不能将道歉减少到零，但是你可以根据自己的需求找到最佳的折中方案：选择一个合适点使得既不能有太多不一致，也不能出现太多可用性问题。

## 信任，但要确认

我们所有关于正确性、完整性和容错性的讨论都是基于某些事情会出错，而其他事情不会出错的假定基础上。我们将这些假设称为我们的系统模型（请参阅第8章“理论系统模型与现实”）：例如，我们应该假设进程可能会崩溃，机器可能突然断电，网络可能会出现不受控制的延迟或丢弃消息。我们也可以假设写入磁盘的数据在强制同步命令`fsync`执行后不会丢失，内存中的数据没有被破坏，并且CPU的乘法指令总是返回正确的结果。

这些假设是相当合理的，因为大部分时间实际情况就是这样。如果我们不得不每秒都在担心计算机会出错，最终将无法完成任何工作。传统上，系统模型采用二元方法处理故障：假定有些事情可能发生，而其他事情是不会发生的。实际上，这更像是一个概率大小问题。关键问题是，违反我们假设的事情是否的确经常在现实中发生。

我们已经看到，即使未发生访问，磁盘上的数据也可能损坏（请参阅第7章“复制与持久性”），并且网络上的数据损坏有时可能会影响TCP校验值（请参阅第8章的“弱的谎言形式”）。也许这才是我们应该更加关注的事情？

我过去曾经使用的一个应用程序收集了来自客户端的崩溃报告，其中一些报告只能通过设备内存发生了随机位翻转来解释。这听起来难以置信，但是如果有足够多的设备来运行你的软件，即使再不可能的事情也有发生的可能。除了由于硬件故障或辐射导致的随机内存损坏之外，某些不合规的存储器访问模式可以在内存完好情况下导致位翻转<sup>[62]</sup>，这种效应可以用来破坏操作系统中的安全机制<sup>[63]</sup>（称为`rowhammer`）。所以如果深入观察，硬件本身并不是它看起来那样完美正确的抽象。

要说明的是，随机位翻转在现代硬件上的确非常少见<sup>[64]</sup>。我只是想指出，它们是可能会发生的，所以应该值得注意。

## 软件缺陷时的完整性

除了硬件问题，总是存在软件错误的风险（bug），这些错误无法通过底层的网络、



内存或文件系统校验和来捕获。即使是广泛使用的数据库软件也存在bug，例如尽管MySQL和PostgreSQL这种非常健壮且广为人知的数据库，多年来经受了很验证，我个人还是曾经看到过MySQL没有正确维护唯一性约束<sup>[65]</sup>和PostgreSQL的串行化隔离级别下所出现的写倾斜异常<sup>[66]</sup>。而对于那些不太成熟的软件中，情况可能会更糟。

即使在软件设计、测试和代码检查方面投入相当多的努力，但bug仍无法彻底避免。虽然可能性低，且最终可以被发现然后修复，但此前总存在一段时间，可能会造成数据破坏。

当谈及应用程序代码时，我们往往不得不假定里面存在更多的bug，主要是大多数应用程序所做的代码检查和测试的量级无法和数据库相提并论。许多应用程序甚至不能正确使用数据库提供的完整性相关的功能，比如外键或唯一性约束<sup>[36]</sup>。

ACID意义上的一致性（请参见第7章“一致性”）假定数据库以一致状态启动，而事务将其从某一个一致状态转换为另一个一致状态。所以，我们期望数据库始终处于某种一致状态。然而，前提是事务不存在bug。如果应用程序以一种错误地方式使用数据库，例如不安全地使用弱级别隔离，就不能保证数据库的完整性。

## 不要盲目信任承诺

硬件和软件并不能总是处于理想状态，数据损坏迟似乎只是迟早的事情而无法避免。因此，我们至少需要有办法来查明数据是否已经损坏，以便之后修复这些数据，并试图找出错误的根源。检查数据的完整性也被称为审计。

正如第11章“不可变事件的优势”中所描述的，审计不仅适用于财务应用程序。但是，可审计性对于金融来说尤其重要，因为每个人都清楚总会有错误，而且都意识到检测和解决问题的必要性。

成熟的系统同样会考虑不太可能的事情出错的可能性，并且主动管理这种风险。例如，HDFS和Amazon S3等大型存储系统不完全信任磁盘：它们运行后台进程，不断读取文件，将其与其他副本进行比较，并将文件从一个磁盘移动到另一个磁盘，以减轻无提示数据损坏的风险<sup>[67]</sup>。

如果想确保你的数据仍然在那，只能不断地去读取和检查。大多数情况下，情况一起正常，但万一发现异常，则越早发现问题越好。基于此，今早尝试从备份来恢复数据，否则当你丢失数据，你会发现连备份也已经破坏，那时将为时已晚。千万不要盲目地相信系统总是正常工作。





## 验证的文化

像HDFS和AWS S3这样的系统也必须假定磁盘大部分时间都能正常工作，这是一个合理的假设，但这不同于假设它们总是正常工作。然而，目前还没有多少系统秉持这种“信任，但要确认”的理念来不断自我审计。许多人认为正确性的保证是绝对的，而没有为少见但可能的数据损坏而有所准备。我希望将来会有更多的自我验证或自我审计系统（*self-validating*或*self-auditing*），不断的检查自身的完整性，而不是依赖盲目的信任<sup>[68]</sup>。

我担心推崇ACID数据库的文化会导致我们在盲目信任底层技术（如事务机制）的基础上开发应用程序，而忽视了整个过程中的可审计性。由于我们所信任的这些技术在大多数情况下工作得很好，审计机制似乎并不值得投入。

但之后数据库的格局发生了显著变化：在NoSQL下，弱一致性成为新常态；底层存储技术并非那么成熟可靠但却普遍使用。然而，数据审计机制尚未形成，我们仍然继续盲目信任底层技术而开发上层应用，这种做法已经变得更加危险。我们应该花些时间来思考一下关于可审计性的设计。

## 可审计性的设计

如果某个事务改变了多个数据库对象，就很难说清楚事务究竟意味着什么了。即使捕获了事务日志（请参阅第11章的“变更数据捕获”），相关表的插入、更新和删除操作等并不一定能清楚地说明为什么这样执行。决定这些变化的应用逻辑的调用总是千变万化，无法轻易复现。

相比之下，基于事件的系统可以提供更好的可审计性。在事件源方法中，用户对系统中的输入都被表示为一个单一的不可变事件，并且任何结果状态的更新都是依据该事件派生而来。派生可以很确定性的执行并且是可重复的，所以通过相同版本的派生代码来处理相同的事件日志将产生相同的状态更新。

清楚地控制数据流（请参阅第10章“批处理输出的哲学”）可以使数据的来源管理更加清晰，从而使完整性检查更加可行。对于事件日志，我们可以使用哈希校验来检查存储层是否发生数据破坏。对于派生状态，我们可以重新运行对相同的事件日志执行的批处理和流处理，以检查是否得到相同的结果，甚至是并行运行一个冗余派生系统。

确定的和定义清晰数据流也有助于系统调试和跟踪系统的操作，从而确定为什么发生了某些事情<sup>[4,69]</sup>。如果中间发生了意外事件，可以提供诊断能力来重现导致意外事件的相同环境，这种精准复现历史时刻的调试能力将非常有价值。





## 端到端论点的再讨论

如果我们不能完全相信系统中的每个组件都能免于损坏（每一个硬件永不出错，且每一个软件都没有bug），那么我们至少也要定期检查数据的完整性。如果不检查的话，就无法发现数据破坏直到为时已晚，并造成了一些相关损失，此时再去追踪这个问题就会更加困难和付出更多代价。

检查数据系统的完整性最好以端到端的方式进行（请参阅本章前面的“数据库端到端的争论”）：在完整性检查中所包含的系统部件越多，则过程中某些阶段发生无告警的数据破坏的概率就越少。如果我们可以检查整个派生系统流水线是端到端正确的，那么路径中的任何磁盘、网络、服务和算法等已经全部囊括在内了。

持续的端到端完整性检查可以提高你对系统正确性的信心，从而使你的发展速度更快<sup>[70]</sup>。与自动化测试一样，审计增加了发现错误的可能性，从而降低了系统变更或新存储技术潜在造成损害的风险。如果你不害怕做出改变，就可以更好地演化应用使之满足不断变化的需求。

## 审计数据系统的工具

目前，将可审计性列为高优先级关注的数据库系统并不多。有些应用程序实现了内部的审计机制，例如将所有更改记录到单独的审计表中，但是保证审计日志的完整性和数据库状态仍然有些困难。事务日志可以通过定期使用硬件安全模块对事务日志进行签名来防止篡改，但这并不能保证正确的事务首先写入日志。

使用加密工具来证明系统的完整性将是很有趣的，这种方式对于广泛的硬件和软件问题甚至是潜在的恶意行为都是可靠的。加密货币，区块链和分布式账本技术，比如比特币、以太坊、波纹、恒星等<sup>[71-73]</sup>已经在这个领域崭露头角。

我没有资格评论这些技术作为货币或者共识合约的优缺点。但是，从数据系统的角度来看，它们包含了一些有趣的想法。本质上，它们都是分布式数据库，具有数据模型和事务机制，不同的副本可以由不信任的组织托管。副本不断检查彼此的完整性，并采用共识协议来就执行的事务达成一致。

我对这些技术中涉及的拜占庭式容错方面存有怀疑（参见第8章“拜占庭式故障”），而且我认为发现工作证明（比如比特币挖矿）机制非常浪费资源。比特币的交易吞吐量是相当低的，尽管更多是处于出于政治和经济方面的原因。但是，其完整性检查方面很有意思。

密码审计和完整性检查通常依赖于Merkle树<sup>[74]</sup>，这是一种哈希散列树，可以用来有效





地证明某个数据集（和其他一些数据集）中出现的记录。除了加密货币的炒作之外，证书透明度是一种依赖Merkle树来检查TLS/SSL证书有效性的安全技术<sup>[75,76]</sup>。

我可以想象完整性检查和审计算法，如证书透明度和分布式账本，在通用的数据系统中将得到越来越广泛的应用。使它们具有和没有加密审计的系统同样的可扩展性还需要更多一些工作，并且还要保持尽可能低的性能损失。但我认为这是一个值得关注的领域。

## 做正确的事情

在本书的最后部分，我想退后一步再谈谈。本书中，我们已经考察了数据系统的各种不同架构，评估其优缺点，探讨了如何构建可靠、可扩展和可维护的应用系统。但是，我们特意留了一个非常重要的基础部分，现在，在结束之前我想重点讨论。

每个系统的都有其构建目的，我们所采取的每一个行动都会产生有意或无意的后果。目的可能像赚钱一样简单，但对世界带来的影响可能远远超出我们的初衷。建立这些系统的工程师有责任仔细考虑这些后果，并有意识地决定我们想要生活在什么样的世界。

我们将数据作为一个抽象的东西来讨论，但要记住，许多数据集都是关于人的：他们的行为、他们的兴趣和他们的身份。我们必须以人性和尊重来对待这些数据。用户也是人，人的尊严是最重要的。

软件开发越来越牵扯到重要的道德选择。有一些指导方针可以帮助软件工程师解决问题，比如ACM的软件工程道德规范和专业实践<sup>[77]</sup>，但是在实践中却很少被讨论、应用和实施。因此，工程师和产品经理们有时对隐私和产品潜在的负面后果采取了非常漫不经心的态度<sup>[78-80]</sup>。

技术不合适或者技术本身有缺陷并不重要，重要的是如何使用它以及如何影响人们。对于像搜索引擎这样的软件系统来说，就像枪支武器一样重要。我认为软件工程师如果只专注于技术而忽视其后果是不够的，道德责任也是我们要担起的责任。评判道德总是困难的，但它太重要了以至无论如何不能被忽视。

## 预测性分析

预测分析是“大数据”炒作的主要部分。使用数据分析预测天气或疾病的传播是一回事<sup>[81]</sup>，预测一个囚犯是否可能再次犯罪，贷款申请人是否有可能违约，或者一个保险客户是否有可能进行昂贵的索赔，则是另一回事。后者将直接影响到个人的生活。





在线支付网络需要防止欺诈性交易，银行希望避免不良贷款，航空公司希望避免劫持，公司希望避免雇用无用或不值得信任的人，这些都是很自然的事情，从它们的角度来看，错失商机的代价很低，但坏账或者有问题员工的成本要高得多，所以相关组织谨慎起见是很自然的。如果有疑问，他们最好说不。

然而，基于算法的决策变得越来越普遍，被这些算法（准确地或错误地）标记为有风险的人可能会由此面临大量“不”的遭遇。系统地被排除在包括工作、航空旅行、保险、物业租赁、金融服务以及其他社会关键方面之外，会对一个人的自由生活产生巨大的制约，因此也被称为“算法囚笼”<sup>[82]</sup>。在很多国家，司法制度会在证明有罪之前推定当事人无罪。然而，自动化的系统则有可能系统地、任意地排除某个人参与社会活动，而且是在这个人没有任何犯罪证据的情况下，并且对他/她来说几乎没有上诉的机会。

## 偏见与歧视

算法锁做出的决定不一定比人类做得更好或更糟。每个人都可能有偏见，即使他们积极地试图消除自己的偏见，歧视性做法也可能在文化上形成制度化。希望根据数据而不是主观和本能的评估来做出决定可能会更加公平，而且给了那些在传统体制中经常被忽视的人以更好的机会<sup>[83]</sup>。

当开发预测分析系统时，我们不仅仅是通过使用软件来指定什么时候说是或什么时候说否的规则来自动化决策，甚至将规则本身从数据中推断出来。但是，这些系统所学到的模式是不透明的：即使数据中存在一些相关性，我们也可能不知道为什么。如果在算法的输入中存在系统性偏见，那么系统很可能吸收塔并在最终输出中放大这种偏见<sup>[84]</sup>。

在许多国家，反歧视法律禁止根据种族、年龄、性别、性取向、残疾或信仰等受保护的点而区别对待不同的人。可以分析一个人的数据的其他特征，但是如果他们与受保护的点相关，会发生什么？例如，在某些敏感的地区，一个人的邮政编码，甚至他们的IP地址，都是一个强有力的预测。就像这样，相信一个算法可以以某种方式将偏倚的数据作为输入并产生公平和公正的输出<sup>[85]</sup>，这似乎是荒谬的。然而，这种观点似乎常常被数据驱动型决策的支持者所采纳，这种态度被讽刺为“机器学习就像为偏见在洗钱”<sup>[86]</sup>。

预测分析系统只是基于过去而推断，如果过去是有偏见的，它们就会把这种偏见编码下来。如果我们希望未来比过去更好，那么就需要道德想象力，而这只有人类才具备<sup>[87]</sup>。数据和模型只应该是我们的工具，而不是我们的主人。





## 责任与问责

自动决策引发了责任与问责方面的问题<sup>[87]</sup>。如果一个人犯了错误，他们可以被追究责任，受决定影响的人可以上诉。算法也会犯错误，但是如果出了问题谁会负责<sup>[88]</sup>？当一辆自动驾驶汽车引发事故时，谁来负责？如果自动信用评分算法系统地区分特定类型的人，是否有任何追索权？如果你的机器学习系统的决定受到司法审查，你能向法官解释算法是如何做出决定的吗？

信用评级机构是收集数据并基于数据做出对人相关决策。不良的信用评分会让生活变得困难，但至少一个信用评分通常是基于一个人的实际借款历史等相关事实，并且记录中的任何错误都可以被纠正（虽然这些机构通常不会让这种修改变得那么容易）。然而，基于机器学习的评分算法通常使用更广泛的输入范围，而且更加不透明，更难理解某个特定决策是如何发生的，以及是否有人受到不公正的对待<sup>[89]</sup>。

信用评分总结了“你以前的行为”，而预测分析则通常基于“谁与你类似，以及和你相似的人在过去的行为如何”。与他人行为相似的对比意味着刻板的来区分人们，例如基于他们居住的地方（这是一个接近社会经济阶层的信息）。那么那些被放错位置的人呢？而且，如果由于错误的数据而做出的决定是不正确的，则追索权几乎是不可能的<sup>[87]</sup>。

很多数据本质上是统计出来的，这意味着即使概率分布总体上是正确的，个别情况也可能是错的。例如，如果某个国家的平均寿命是80岁，并不意味着某个人过了80岁生日就会死亡。基于平均分布和概率分布无法预测某个人的寿命。同理，预测系统的输出是概率性的，在个别情况下也可能是错误的。

盲目地相信数据至高无上不仅是误解的，而且是非常危险的。随着数据驱动的决策变得越来越普遍，我们需要弄清楚如何使算法更负责任和透明，如何避免强化现有的偏见，以及如何在错误不可避免时加以修复。

我们还需要弄清楚如何防止数据被滥用，并努力发挥数据的正面作用。例如，通过分析可以揭示人们生活中的财务状况和社会特征。一方面，可以利用它把援助和支持集中去帮助那些最需要援助的人。另一方面，它有时被不正当的企业用来识别弱势群体，并向他们销售高风险的产品，比如高成本的借贷和毫无价值的大学学位<sup>[87,90]</sup>。

## 反馈环路

即使预测性应用对人们所产生的立竿见影的长远影响较少，比如推荐系统，我们仍然要面对一些难以解决的问题。如果预测用户希望看到什么的服务做得很好，他们最终可能只会向人们展示系统认同的观点，从而容易产生陈旧观念、错误信息和两极分化



的回声室效应。我们已经看到社交媒体在竞选活动上产生的类似影响<sup>[91]</sup>。

当预测分析影响人们的生活时，特别是由于自我强化反馈环路而出现一些有害问题。比如，考虑雇主使用信用评分来评估潜在雇员。你可能是一个有良好信用评分的员工，但是由于一些不可控情况，你会突然发现自己陷入了经济困境。由于错过了信用卡账单还款，你的信用评分会受损，并可能因此找工作会变得困难。失业使你陷入贫困，使信用评分进一步恶化，结果就是更难找到工作<sup>[87]</sup>。由于不合适的假设，产生了这样一个隐藏在数学严谨性和数据的伪装背后的下降旋涡。

我们不是总能预测什么时候发生这样的反馈环路。然而，通过思考整个系统（不仅是计算机化的部分，还有与之互动的人）可以预测许多后果，这是一种被称为系统思维的方法<sup>[92]</sup>。我们可以尝试理解一个数据分析系统是如何响应不同的行为、结构和特征。系统是否强化和扩大了人们之间存在的差异（例如，使富者变得更富或穷人变得更穷）？还是试图打击不公平性？即使有最好的意图，我们也必须小心意外的后果。

## 数据隐私与追踪

除了预测分析方面的问题（即使用数据来做出关于人的自动化决策）之外，数据收集本身也存在道德问题。收集数据的组织与正在收集数据的人之间的是什么关系呢？

当系统只保存用户所明确输入的数据（因为他们希望系统以某种方式进行存储和处理这些数据），系统就是在为用户执行一项服务：用户就是客户。但是，处于其他目的而对用户的活动被跟踪并记录时，这种关系就不那么清晰了。服务不再仅仅是用户告诉它所做的事情，而是服务于自身的利益，这可能与用户的利益相冲突。

跟踪行为数据对于许多面向用户的在线服务变得越来越重要：跟踪哪些搜索结果被单击有助于提高搜索结果的排名；推荐“喜欢X的人也喜欢Y”，可以帮助用户发现有趣而有用的东西；A/B测试和用户流量分析可以帮助指出如何改进用户界面。这些功能需要一定量的用户行为跟踪，并且用户也可以从中受益。

但是，根据公司的商业模式，追踪往往不止于此。如果服务是通过广告获得资助的，那么广告主就是实际的客户，而用户的利益则是次要的。跟踪数据会更加详细，分析变得更加深入，数据也会被保留很长一段时间，以便为营销目标去建立每个人的详细资料。

现在，公司和被收集数据的用户之间的关系开始变得和以往大不一样了。用户得到免费的服务，并尽可能地被引诱参与到服务中。对用户的追踪不再是服务与个人，而是服务于资助广告客户的需求。我认为这种关系可以用一个更阴暗的词来描述：监视。



## 监控

我们做一个有意思的实验，尝试用监控一词来代替数据，观察常见的短语是否听起来还是那么美好<sup>[93]</sup>。比如说：“在我们的监控驱动的组织，我们收集实时监控流，并将其存储在我們的监控仓库。我们的监测科学家使用先进的分析和监测处理，以获得新的见解。”

把这个观念实验用于本书书名“数据密集型应用系统设计”同样会引发巨大争议，但是我认为替换一个更强力的词汇更能说明其中的道理。在我们试图依靠软件“掌控世界”<sup>[94]</sup>过程中，我们建立了世界上迄今为止最大规模的监视基础设施。我们正朝着物联网快速迈进，每个有人居住的空间至少包含一个联网的麦克风，很多设备包括智能手机，智能电视，语音辅助设备，婴儿监视器，甚至是儿童玩具等都在使用基于云的语音识别。这些设备中的很多都有可怕的安全记录<sup>[95]</sup>。

即使是条件受限制的地方，也梦想着在每个房间放置一个麦克风，并强迫每个人不断地携带一个能跟踪他们位置和动作的设备。然而，我们今天显然自愿地、甚至迫不及待地把自己置身于一个全面监视的世界里。不同之处在于数据是由很多公司收集而不是政府机构<sup>[96]</sup>。

并非所有的数据收集都必须满足监控的要求，但是检查这些数据可以帮助理解我们与数据收集者的关系。为什么我们看起来乐意接受企业的监督？也许你觉得没有什么可隐瞒的东西，换句话说，你完全和现有的组织结构是一个阶层的，你不是被边缘化的少数派，不必害怕<sup>[97]</sup>。但是不是每个人都如此幸运。也许这是因为你知道监控的目的似乎是好的，不是强制性和必须顺从的，而只是更好的建议和更个性化的营销。但是，结合上一节的预测分析的讨论，它们的差别似乎微乎其微。

我们已经看到汽车保险费用已经和汽车内安装的追踪设备连在一起了，以及健康保险的覆盖范围也取决于人们身上佩戴的健身追踪设备。当监控被用来确定生活中重要的事情，例如保险或就业等方面的东西时，它就开始变得不那么亲切了。此外，数据分析可以揭示出令人惊讶的侵入性的事情：比如，智能手表或健身追踪器中的运动传感器可以相当准确地计算出正在输入的内容（例如密码）<sup>[98]</sup>。况且现在的分析算法变得越来越强大。

## 赞成与选择的自由

我们可能会宣称，用户自愿选择使用跟踪活动的服务，并且已经同意服务条款和隐私政策，因此他们同意被收集数据。我们甚至可以声称，用户正在用他们提供的数据来换取有价值的服务，而为了正常提供服务，跟踪是必要的。毫无疑问，社交网络、搜





索引引擎以及其他各种免费的在线服务对于用户来说的确是有价值，但是即使这样，上述说法也存在问题。

用户几乎不知道什么样的个人数据会进入到数据库，或者数据是如何保留和处理的，大多数隐私政策的条款也极尽所能地搞得含混不清。不清楚他们的数据会发生什么，用户就不能给予任何有意义的认同。通常，来自用户的数据还被用到了不是该服务的用户身上，并且该用户根本就没有同意数据收集的任何条款。我们在本书中讨论的派生数据集（其中来自整个用户群的数据可能是行为跟踪和外部数据源的结合体）恰恰是用户无法获得任何有意义的理解的那种。

而且，数据是通过单向过程从用户提取而来，而不是通过真正的互惠关系，也不是公平的价值交换。没有对话，用户无法选择提供多少数据以及他们会收到什么样的服务：服务与用户之间的关系是非常不对称的：这些条款是由服务提供商所设置，而不是由用户<sup>[99]</sup>。

对于不同意被监视的用户，唯一真正的选择就是不使用服务。但是这个选择也不是免费的：如果一项服务非常受欢迎以至于“被大多数人认为是基本的社会参与所需要的”<sup>[99]</sup>，那么指望人们选择退出这项服务是不合理的，使用它变成一种事实性强制约束。例如，在大多数社会群体中，携带智能手机，使用Facebook进行社交以及使用Google查找信息已成为常态。特别是当一个服务具有网络效应时，人们选择不使用它是有社会成本的。

由于担心服务跟踪用户而决定拒绝使用，这只对极少数拥有足够的时间和知识来充分了解隐私政策的人群可以称得上是一种选择，并且他们可以不需要担心由此可能会失去某些机会而被迫参与这些服务。然而，对于处境较差的人来说，选择自由没有意义：对他们来说，被监视变得不可避免。

## 数据隐私和使用

有时人们会说“再无隐私”，理由是一些用户愿意将他们生活的各种事情发布到社交媒体上，有时是很稀松平常的事情，有时甚至是非常私密的。但我认为，这种说法是错误的，因为它基于了对隐私这个词的误解。

拥有隐私并不意味着一切事情都要保密；它意味着你可以自由选择向谁展示，并展示哪些东西，要公开什么，以及要保密什么。隐私权是一个决定权：每个人都能够决定在各种情况下如何在保密和透明之间取舍<sup>[99]</sup>。这事关个人的自由和自主。

当通过监控基础设施从人们身上提取数据时，隐私权不一定被破坏，而是转移到了



数据收集者。获取数据的公司基本上都会说：“相信我们对你的数据只做正确的事情”，这意味着决定要披露什么以及保密的权利是从个人转移到公司了。

这些公司最终选择对大部分数据继续保持私密，因为泄露数据会引起可怕的后果，并且会损害它们的商业模式（这样它们就可以比其他公司更了解用户）。关于用户的隐私信息通常是间接地被泄露，例如借助数据分析，将广告投放给特定人群（例如患有特定疾病的人）。

即使针对特定人群的广告并没有识别出特定的个体，但他们也丧失了披露隐私信息的权利，例如是否患有某种疾病。现在不是用户根据自己的喜好决定向谁显示，显示哪些隐私内容；做决定的是公司，公司通常会以最大化利润为目标来处置隐私权。

许多公司都不希望人们将他们视为作恶者，例如想法避免被问及那些具有侵犯性的数据收集，而让人们感觉他们是服务于用户体验或感受。即使这些体验有时很糟糕：例如，有些事情虽然是事实，但它触发人们内心痛苦的回忆，用户可能并不希望它被反复提及<sup>[100]</sup>。对于任何类型的数据，我们都应该期望在某种程度上会出现错误、不受欢迎或不适当的可能，并建立响应的处理机制。当然，何为“不受欢迎”或“不适当”，完全是由人来判断决定的。除非我们明确地设计程序使得它们尊重人类的需求，否则算法对此是没有任何概念的。作为这些系统的构建者，我们必须谦虚、乐于接受并为之做好准备。

所谓隐私设置，即允许在线服务的用户控制哪些数据对他人可见，只是将一些控制权还给用户的起点。但是，无论设置如何，服务本身仍然可以不受限制地访问数据，并且可以以隐私策略所允许的任何方式自由地使用它。即使服务承诺不会将数据出售给第三方，它通常也授予自身无限制的权利，在内部处理和分析数据时往往远远超过用户公开可见的权利。

这种大规模地将隐私控制权从个人转移到公司在历史上是绝无仅有的<sup>[99]</sup>。监视一直存在，但过去是昂贵和手动的，不可扩展和自动化。信任关系一直存在，例如病人与医生之间，或被告与律师之间，但在这些情况下，数据的使用受到道德、法律和监管限制的严格管制。互联网服务使得在没有用户同意的情况下积累大量敏感信息更加容易，并且在用户不知情相关后果的前提下大规模地使用它。

## 数据作为资产和权力

由于行为数据是用户与服务交互作用过程中产生的副产品，有时称为“废弃数据”，表示这类数据是无价值的。通过这种方式，行为分析和预测分析可以被看作是某种回收方式，从那些本来可能被抛弃的数据中提取价值。



更准确的说是应该反过来看：从经济的角度来看，如果有针对性的广告是一种服务，那么关于人的行为数据就是服务的核心资产。在这种情况下，与用户交互的应用程序仅仅是一种诱使用户将越来越多的个人信息提供给监控基础设施的手段<sup>[99]</sup>。在线服务中经常可以看到令人愉悦的人类创造力和社交关系被数据提取机器不负责任地利用着。

数据中介公司的存在也印证了个人数据是宝贵资产的说法，这个数据中间商是一个秘密行业，从事采购、汇总、分析、推断和兜售侵入性个人数据，主要是为了营销目的<sup>[90]</sup>。很多初创公司主要靠它们的用户量（或者“眼球”，即他们的监视渗透能力）来估价。

因为这些数据是有价值的，所以很多人都想要它。公司很明显是需要这些数据的，这就是为什么他们收集它的原因。当公司破产时，收集到的个人资料就是被出售的资产之一。而且，数据很难保证安全，所以令人不安的破坏事件时常发生<sup>[102]</sup>。

这些观察使得某些批评者认为数据不仅仅是一种资产，而是一种“有毒资产”<sup>[101]</sup>，或者至少是“有害物质”<sup>[103]</sup>。即使我们认为有能力防止数据滥用，但是每当我们收集数据时，都需要平衡带来的好处和落入坏人手中的风险：计算机系统可能会被犯罪分子篡改，数据可能会被内部人员泄露，公司可能会落入不当价值观的无良管理层手中等。

收集数据时，一定要综合考量<sup>[104]</sup>。

正如古老的格言所言，“知识就是力量”。此外，“审视他人但避免自我审查是最重要的权力形式之一”<sup>[105]</sup>。尽管今天的科技公司并没有公开地寻求某些权力，但是它们所积累的数据和知识给了它们很大的控制权力，而且很多是在私下进行，不在公众监督之内。

## 记住工业革命

数据是信息时代的关键性特征。互联网，数据存储，处理器和软件驱动的自动化正在对全球经济和人类社会产生重大影响。考虑到过去十余年我们的日常生活和社会组织所发生的巨大变化，有理由相信在未来几十年可能会继续发生根本性的变化，由此不由得联想到工业革命<sup>[87,96]</sup>。

工业革命是伴随着重大技术和农业经济的大发展而来的，长远来看人们生活水平明显提高。然而，它也带来了一些重大问题：空气污染（由于烟雾和化学过程）以及水污染（由于工业废物和生活废物等）的后果是可怕的。工厂老板光鲜亮丽，工人住房条





件差、工作时间长、工作条件恶劣。不合法的工人到处可见，包括在矿场的危险的，低薪的工作。

人类花了很长时间才制定了一系列的保护措施，如环境保护条例、工作场所安全协议、禁止使用童工和食品卫生检查规范等。毫无疑问，当工厂不能将废物倾倒入河中，不能销售受污染的食品，也不能再继续剥削工人时，成本就会增加。但是，整个社会获得了巨大的受益，所谓“绿水青山就是金山银山”，几乎没有人再愿意回到之前污染的年代<sup>[87]</sup>。

正如工业革命存在需要被管理的黑暗面一样，向信息时代的过渡也有需要面对和解决的重大问题。我相信收集和使用数据就是其中一个。用Bruce Schneier<sup>[96]</sup>的话来说：

数据问题是信息时代的污染问题，保护隐私类似包含环境。几乎所有的电脑都在产生信息，数据持续不断然后又慢慢消逝。我们如何定义这个问题以及如何处置对信息经济的健康运行至关重要。正如我们今天回顾工业时代的早期几十年，惊呼为什么我们的祖先如此急于建立工业世界而忽视了污染问题，我们的子孙也将审视现在的信息时代初期，并评判我们如何来应对数据收集与滥用。

我们理应做到让后代们感到骄傲。

## 立法与自律

数据保护法可能有助于维护个人的权利。例如，1995年的“欧洲数据保护指令”规定，个人数据必须“为特定的，明确的和合法的目的而收集，而不能与这些目的不相符的方式来处理”，而且数据必须“就其收集的目的而言充分，相关但不过分”<sup>[107]</sup>。

但是，这个立法在今天的互联网环境下是否有效还是值得商榷<sup>[108]</sup>。这些规则直接违背了大数据的理念，即最大化数据收集，将其与其他数据集合在一起，进行实验和探索，以产生新的见解。探索意味着将数据用于不可预见的目的，这与用户同意的“明确的和清晰的”目标相反（如果我们充分理解所认同内容<sup>[109]</sup>）。目前更新版的规定仍在制定之中<sup>[89]</sup>。

收集大量数据的公司反对监管，它们认为这是对创新的负担和阻碍。在某种程度上，这种反对是有道理的。例如，分享医疗数据时，隐私风险明显，但也有潜在的机会：如果数据分析能够帮助我们实现更好的诊断或找到更好的治疗方法，可以预防多少人死亡<sup>[110]</sup>？过度管制可能会阻止这种突破。这种潜在的机会与风险之间平衡难以很好掌握<sup>[105]</sup>。

从根本上说，我认为需要对针对个人数据的技术领域有观念上转变。我们应该停止过



度以用户为衡量指标，牢记用户值得尊重。我们应该主动调整数据收集和处理流程，建立和维持与那些依赖我们软件的人们之间的信任关系<sup>[111]</sup>。我们应该主动向用户介绍他们的数据如何使用，而不是让他们蒙在鼓里全然不知。

我们应该允许每个人维护自己的隐私，即控制自己的数据而不是通过监视来窃取他们的控制权。用户掌控自己数据的个人权利就像是国家公园中的优美环境：如果对环境没有明确的保护和关心，绿水青山将会被破坏。那将是所有公众的悲哀，我们所有人都会受到影响。无所不在的监视和数据滥用并非不可避免，我们仍然能够阻止它。

我们究竟该如何实现上述目标还是一个开放的问题。首先一点，我们不应该永远保留数据，一旦不再需要，就尽快清除它们<sup>[111,112]</sup>。清除数据与不可变性的想法背道而驰（请参阅第11章“不可变性的限制”），但是可以解决这个问题。我所看到的一个很有前途的方法是通过加密协议来实施访问控制，而不仅仅是通过策略<sup>[113,114]</sup>。总的来说，观念与态度的变化都是必要的。

## 小结

本章，我们讨论了设计数据系统的新方法，以及包括了我个人对未来的看法和猜测。我们从观察开始，即没有一种工具可以有效地服务于所有可能的场景，因此应用程序必须编排多个不同的软件来完成他们的目标。我们讨论了如何使用批处理和事件流来解决这个数据集成问题，以使数据在不同系统之间灵活流动。

基于这种方法，某些系统被指定为记录系统，而其他数据则是通过转换而来。通过这种方式，我们可以维护索引、实体化视图、机器学习模型、统计摘要等。通过使这些推导和变换异步松耦合，防止了局部问题可能扩散到系统的不相关部分，从而提高了整个系统的鲁棒性和容错性。

将数据流表示为从一个数据集到另一个数据集的转换也有助于演化应用程序：如果要更改其中一个处理步骤，例如更改索引或缓存的结构，则可以在整个输入数据集上重新运行新的转换代码以便重新输出。同样，如果出现问题，可以修复代码然后重新处理数据来执行恢复。

这些过程与内部数据库完成的过程非常相似，因此我们将数据流应用系统的概念重新分解为数据库组件，并通过组合这些松耦合的组件来构建应用程序。

派生状态可以通过观察基础数据的变化来不断更新。而且，下游消费者可以进一步观察派生的状态。我们甚至可以将此数据流一直传送到显示数据的最终用户设备，从而构建动态更新的用户界面以反映数据更改并支持离线使用。



接下来，我们讨论了在出现故障时如何确保处理正确。我们看到，强的完整性保证可以通过异步事件处理，通过使用端到端操作标识使操作最终满足幂等性，并通过异步检查约束来实现。客户可以等到检查通过或者不用等，但是如果万一违反约束则需要事后道歉。这种方法比使用分布式事务的传统方法更具可扩展性和可靠性，并且适合于实践中有多个业务流程同时工作的场景。

围绕数据流来构建应用程序并异步检查约束，我们可以避免大多数协调工作，系统保证完整性并性能良好，即使在地理位置分散的情况下或是出现故障时也是如此。然后，我们简单介绍了使用审计来验证数据的完整性并检测数据是否发生破坏。

最后，我们又退后一步，审视了构建数据密集型应用系统在道德层面的一些问题。我们看到，虽然数据可以用来帮助人们，但是也可能造成重大的伤害：做出严重影响人们生活的貌似公平的决定，这种算法决定难以对其提起诉讼；导致歧视和剥削；使监视泛滥；暴露私密信息等。我们也面临着数据泄露的风险，而且即使善意的数据使用也可能会产生某些意想不到的后果。

由于软件和数据对世界的影响如此之大，我们的工程师们必须谨记，我们有责任为我们赖以生存的世界而努力：一个以人性和尊重来对待人的世界。我希望我们能够一起为实现这一目标而努力。

## 参考文献

- [1] Rachid Belaid: “Postgres Full-Text Search is Good Enough!,” *rachbelaid.com*, July 13, 2015.
- [2] Philippe Ajoux, Nathan Bronson, Sanjeev Kumar, et al.: “Challenges to Adopting Stronger Consistency at Scale,” at *15th USENIX Workshop on Hot Topics in Operating Systems (HotOS)*, May 2015.
- [3] Pat Helland and Dave Campbell: “Building on Quicksand,” at *4th Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2009.
- [4] Jessica Kerr: “Provenance and Causality in Distributed Systems,” *blog.jessitron.com*, September 25, 2016.
- [5] Kostas Tzoumas: “Batch Is a Special Case of Streaming,” *data-artisans.com*, September 15, 2015.





- [6] Shinji Kim and Robert Blafford: “Stream Windowing Performance Analysis: Concord and Spark Streaming,” *concord.io*, July 6, 2016.
- [7] Jay Kreps: “The Log: What Every Software Engineer Should Know About Real-Time Data’s Unifying Abstraction,” *engineering.linkedin.com*, December 16, 2013.
- [8] Pat Helland: “Life Beyond Distributed Transactions: An Apostate’s Opinion,” at *3rd Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2007.
- [9] “Great Western Railway (1835-1948),” Network Rail Virtual Archive, *networkrail.co.uk*.
- [10] Jacqueline Xu: “Online Migrations at Scale,” *stripe.com*, February 2, 2017.
- [11] Molly Bartlett Dishman and Martin Fowler: “Agile Architecture,” at *O’Reilly Software Architecture Conference*, March 2015.
- [12] Nathan Marz and James Warren: *Big Data: Principles and Best Practices of Scalable Real-Time Data Systems*. Manning, 2015. ISBN: 978-1-617-29034-3.
- [13] Oscar Boykin, Sam Ritchie, Ian O’Connell, and Jimmy Lin: “Summingbird: A Framework for Integrating Batch and Online MapReduce Computations,” at *40th International Conference on Very Large Data Bases (VLDB)*, September 2014.
- [14] Jay Kreps: “Questioning the Lambda Architecture,” *oreilly.com*, July 2, 2014.
- [15] Raul Castro Fernandez, Peter Pietzuch, Jay Kreps, et al.: “Liquid: Unifying Nearline and Offline Big Data Integration,” at *7th Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2015.
- [16] Dennis M. Ritchie and Ken Thompson: “The UNIX Time-Sharing System,” *Communications of the ACM*, volume 17, number 7, pages 365-375, July 1974. doi: 10.1145/361011.361061.
- [17] Eric A. Brewer and Joseph M. Hellerstein: “CS262a: Advanced Topics in Computer Systems,” lecture notes, University of California, Berkeley, *cs.berkeley.edu*, August 2011.
- [18] Michael Stonebraker: “The Case for Polystores,” *wp.sigmod.org*, July 13, 2015.



- [19] Jennie Duggan, Aaron J. Elmore, Michael Stonebraker, et al.: “The BigDAWG Polystore System,” *ACM SIGMOD Record*, volume 44, number 2, pages 11-16, June 2015. doi:10.1145/2814710.2814713.
- [20] Patrycja Dybka: “Foreign Data Wrappers for PostgreSQL,” *vertabelo.com*, March 24, 2015.
- [21] David B. Lomet, Alan Fekete, Gerhard Weikum, and Mike Zwilling: “Unbundling Transaction Services in the Cloud,” at *4th Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2009.
- [22] Martin Kleppmann and Jay Kreps: “Kafka, Samza and the Unix Philosophy of Distributed Data,” *IEEE Data Engineering Bulletin*, volume 38, number 4, pages 4-14, December 2015.
- [23] John Hugg: “Winning Now and in the Future: Where VoltDB Shines,” *voltldb.com*, March 23, 2016.
- [24] Frank McSherry, Derek G. Murray, Rebecca Isaacs, and Michael Isard: “Differential Dataflow,” at *6th Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2013.
- [25] Derek G Murray, Frank McSherry, Rebecca Isaacs, et al.: “Naiad: A Timely Dataflow System,” at *24th ACM Symposium on Operating Systems Principles (SOSP)*, pages 439-455, November 2013. doi:10.1145/2517349.2522738.
- [26] Gwen Shapira: “We have a bunch of customers who are implementing ‘database inside-out’ concept and they all ask ‘is anyone else doing it? are we crazy?’ ” *twitter.com*, July 28, 2016.
- [27] Martin Kleppmann: “Turning the Database Inside-out with Apache Samza,” at *Strange Loop*, September 2014.
- [28] Peter Van Roy and Seif Haridi: *Concepts, Techniques, and Models of Computer Programming*. MIT Press, 2004. ISBN: 978-0-262-22069-9.
- [29] “Juttle Documentation,” *juttle.github.io*, 2016.
- [30] Evan Czaplicki and Stephen Chong: “Asynchronous Functional Reactive



Programming for GUIs,” at *34th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, June 2013. doi:10.1145/2491956.2462161.

[31] Engineer Bainomugisha, Andoni Lombide Carreton, Tom van Cutsem, Stijn Mostinckx, and Wolfgang de Meuter: “A Survey on Reactive Programming,” *ACM Computing Surveys*, volume 45, number 4, pages 1-34, August 2013. doi:10.1145/2501654.2501666.

[32] Peter Alvaro, Neil Conway, Joseph M. Hellerstein, and William R. Marczak: “Consistency Analysis in Bloom: A CALM and Collected Approach,” at *5th Biennial Conference on Innovative Data Systems Research (CIDR)*, January 2011.

[33] Felienne Hermans: “Spreadsheets Are Code,” at *Code Mesh*, November 2015.

[34] Dan Bricklin and Bob Frankston: “VisiCalc: Information from Its Creators,” *danbricklin.com*.

[35] D. Sculley, Gary Holt, Daniel Golovin, et al.: “Machine Learning: The High-Interest Credit Card of Technical Debt,” at *NIPS Workshop on Software Engineering for Machine Learning (SE4ML)*, December 2014.

[36] Peter Bailis, Alan Fekete, Michael J Franklin, et al.: “Feral Concurrency Control: An Empirical Investigation of Modern Application Integrity,” at *ACM International Conference on Management of Data (SIGMOD)*, June 2015. doi:10.1145/2723372.2737784.

[37] Guy Steele: “Re: Need for Macros (Was Re: Icon),” email to ll1-discuss mailing list, *people.csail.mit.edu*, December 24, 2001.

[38] David Gelernter: “Generative Communication in Linda,” *ACM Transactions on Programming Languages and Systems (TOPLAS)*, volume 7, number 1, pages 80-112, January 1985. doi:10.1145/2363.2433.

[39] Patrick Th. Eugster, Pascal A. Felber, Rachid Guerraoui, and Anne-Marie Kermarrec: “The Many Faces of Publish/Subscribe,” *ACM Computing Surveys*, volume 35, number 2, pages 114-131, June 2003. doi:10.1145/857076.857078.

[40] Ben Stopford: “Microservices in a Streaming World,” at *QCon London*, March 2016.







[41] Christian Posta: “Why Microservices Should Be Event Driven: Autonomy vs Authority,” *blog.christianposta.com*, May 27, 2016.

[42] Alex Feyerke: “Say Hello to Offline First,” *hood.ie*, November 5, 2013.

[43] Sebastian Burckhardt, Daan Leijen, Jonathan Protzenko, and Manuel Fähndrich: “Global Sequence Protocol: A Robust Abstraction for Replicated Shared State,” at *29th European Conference on Object-Oriented Programming (ECOOP)*, July 2015. doi:10.4230/LIPIcs.ECOOP.2015.568.

[44] Mark Soper: “Clearing Up React Data Management Confusion with Flux, Redux, and Relay,” *medium.com*, December 3, 2015.

[45] Eno Thereska, Damian Guy, Michael Noll, and Neha Narkhede: “Unifying Stream Processing and Interactive Queries in Apache Kafka,” *confluent.io*, October 26, 2016.

[46] Frank McSherry: “Dataflow as Database,” *github.com*, July 17, 2016.

[47] Peter Alvaro: “I See What You Mean,” at *Strange Loop*, September 2015.

[48] Nathan Marz: “Trident: A High-Level Abstraction for Realtime Computation,” *blog.twitter.com*, August 2, 2012.

[49] Edi Bice: “Low Latency Web Scale Fraud Prevention with Apache Samza, Kafka and Friends,” at *Merchant Risk Council MRC Vegas Conference*, March 2016.

[50] Charity Majors: “The Accidental DBA,” *charity.wtf*, October 2, 2016.

[51] Arthur J. Bernstein, Philip M. Lewis, and Shiyong Lu: “Semantic Conditions for Correctness at Different Isolation Levels,” at *16th International Conference on Data Engineering (ICDE)*, February 2000. doi:10.1109/ICDE.2000.839387.

[52] Sudhir Jorwekar, Alan Fekete, Krithi Ramamritham, and S. Sudarshan: “Automating the Detection of Snapshot Isolation Anomalies,” at *33rd International Conference on Very Large Data Bases (VLDB)*, September 2007.

[53] Kyle Kingsbury: Jepsen blog post series, *aphyr.com*, 2013-2016.

[54] Michael Jouravlev: “Redirect After Post,” *theserverside.com*, August 1, 2004.

[55] Jerome H. Saltzer, David P. Reed, and David D. Clark: “End-to-End Arguments in



System Design,” *ACM Transactions on Computer Systems*, volume 2, number 4, pages 277-288, November 1984. doi:10.1145/357401.357402.

[56] Peter Bailis, Alan Fekete, Michael J. Franklin, et al.: “Coordination-Avoiding Database Systems,” *Proceedings of the VLDB Endowment*, volume 8, number 3, pages 185-196, November 2014.

[57] Alex Yarmula: “Strong Consistency in Manhattan,” *blog.twitter.com*, March 17, 2016.

[58] Douglas B Terry, Marvin M Theimer, Karin Petersen, et al.: “Managing Update Conflicts in Bayou, a Weakly Connected Replicated Storage System,” at *15th ACM Symposium on Operating Systems Principles (SOSP)*, pages 172-182, December 1995. doi:10.1145/224056.224070.

[59] Jim Gray: “The Transaction Concept: Virtues and Limitations,” at *7th International Conference on Very Large Data Bases (VLDB)*, September 1981.

[60] Hector Garcia-Molina and Kenneth Salem: “Sagas,” at *ACM International Conference on Management of Data (SIGMOD)*, May 1987. doi:10.1145/38713.38742.

[61] Pat Helland: “Memories, Guesses, and Apologies,” *blogs.msdn.com*, May 15, 2007.

[62] Yoongu Kim, Ross Daly, Jeremie Kim, et al.: “Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors,” at *41st Annual International Symposium on Computer Architecture (ISCA)*, June 2014. doi: 10.1145/2678373.2665726.

[63] Mark Seaborn and Thomas Dullien: “Exploiting the DRAM Rowhammer Bug to Gain Kernel Privileges,” *googleprojectzero.blogspot.co.uk*, March 9, 2015.

[64] Jim N. Gray and Catharine van Ingen: “Empirical Measurements of Disk Failure Rates and Error Rates,” Microsoft Research, MSR-TR-2005-166, December 2005.

[65] Annamalai Gurusami and Daniel Price: “Bug #73170: Duplicates in Unique Secondary Index Because of Fix of Bug#68021,” *bugs.mysql.com*, July 2014.

[66] Gary Fredericks: “Postgres Serializability Bug,” *github.com*, September 2015.



- [67] Xiao Chen: “HDFS DataNode Scanners and Disk Checker Explained,” *blog.cloudera.com*, December 20, 2016.
- [68] Jay Kreps: “Getting Real About Distributed System Reliability,” *blog.empathybox.com*, March 19, 2012.
- [69] Martin Fowler: “The LMAX Architecture,” *martinfowler.com*, July 12, 2011.
- [70] Sam Stokes: “Move Fast with Confidence,” *blog.samstokes.co.uk*, July 11, 2016.
- [71] “Sawtooth Lake Documentation,” Intel Corporation, *intelledger.github.io*, 2016.
- [72] Richard Gendal Brown: “Introducing R3 Corda™ A Distributed Ledger Designed for Financial Services,” *gendal.me*, April 5, 2016.
- [73] Trent McConaghy, Rodolphe Marques, Andreas Müller, et al.: “BigchainDB: A Scalable Blockchain Database,” *bigchaindb.com*, June 8, 2016.
- [74] Ralph C. Merkle: “A Digital Signature Based on a Conventional Encryption Function,” at *CRYPTO '87*, August 1987. doi:10.1007/3-540-48184-2\_32.
- [75] Ben Laurie: “Certificate Transparency,” *ACM Queue*, volume 12, number 8, pages 10-19, August 2014. doi:10.1145/2668152.2668154.
- [76] Mark D. Ryan: “Enhanced Certificate Transparency and End-to-End Encrypted Mail,” at *Network and Distributed System Security Symposium (NDSS)*, February 2014. doi:10.14722/ndss.2014.23379.
- [77] “Software Engineering Code of Ethics and Professional Practice,” Association for Computing Machinery, *acm.org*, 1999.
- [78] François Chollet: “Software development is starting to involve important ethical choices,” *twitter.com*, October 30, 2016.
- [79] Igor Perisic: “Making Hard Choices: The Quest for Ethics in Machine Learning,” *engineering.linkedin.com*, November 2016.
- [80] John Naughton: “Algorithm Writers Need a Code of Conduct,” *theguardian.com*, December 6, 2015.







- [81] Logan Kugler: “What Happens When Big Data Blunders?,” *Communications of the ACM*, volume 59, number 6, pages 15-16, June 2016. doi:10.1145/2911975.
- [82] Bill Davidow: “Welcome to Algorithmic Prison,” *theatlantic.com*, February 20, 2014.
- [83] Don Peck: “They’re Watching You at Work,” *theatlantic.com*, December 2013.
- [84] Leigh Alexander: “Is an Algorithm Any Less Racist Than a Human?” *theguardian.com*, August 3, 2016.
- [85] Jesse Emspak: “How a Machine Learns Prejudice,” *scientificamerican.com*, December 29, 2016.
- [86] Maciej Cegłowski: “The Moral Economy of Tech,” *idlewords.com*, June 2016.
- [87] Cathy O’Neil: *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. Crown Publishing, 2016. ISBN: 978-0-553-41881-1.
- [88] Julia Angwin: “Make Algorithms Accountable,” *nytimes.com*, August 1, 2016.
- [89] Bryce Goodman and Seth Flaxman: “European Union Regulations on Algorithmic Decision-Making and a ‘Right to Explanation’,” *arXiv:1606.08813*, August 31, 2016.
- [90] “A Review of the Data Broker Industry: Collection, Use, and Sale of Consumer Data for Marketing Purposes,” Staff Report, *United States Senate Committee on Commerce, Science, and Transportation*, *commerce.senate.gov*, December 2013.
- [91] Olivia Solon: “Facebook’s Failure: Did Fake News and Polarized Politics Get Trump Elected?” *theguardian.com*, November 10, 2016.
- [92] Donella H. Meadows and Diana Wright: *Thinking in Systems: A Primer*. Chelsea Green Publishing, 2008. ISBN: 978-1-603-58055-7.
- [93] Daniel J. Bernstein: “Listening to a ‘big data’ / ‘data science’ talk,” *twitter.com*, May 12, 2015.
- [94] Marc Andreessen: “Why Software Is Eating the World,” *The Wall Street Journal*, 20 August 2011.
- [95] J. M. Porup: “‘Internet of Things’ Security Is Hilariously Broken and Getting





Worse,” *arstechnica.com*, January 23, 2016.

[96] Bruce Schneier: *Data and Goliath: The Hidden Battles to Collect Your Data and Control Your World*. W. W. Norton, 2015. ISBN: 978-0-393-35217-7.

[97] The Grugq: “Nothing to Hide,” *grugq.tumblr.com*, April 15, 2016.

[98] Tony Beltramelli: “Deep-Spying: Spying Using Smartwatch and Deep Learning,” Masters Thesis, IT University of Copenhagen, December 2015. Available at [arxiv.org/abs/1512.05616](https://arxiv.org/abs/1512.05616).

[99] Shoshana Zuboff: “Big Other: Surveillance Capitalism and the Prospects of an Information Civilization,” *Journal of Information Technology*, volume 30, number 1, pages 75-89, April 2015. doi:10.1057/jit.2015.5.

[100] Carina C. Zona: “Consequences of an Insightful Algorithm,” at *GOTO Berlin*, November 2016.

[101] Bruce Schneier: “Data Is a Toxic Asset, So Why Not Throw It Out?,” *schneier.com*, March 1, 2016.

[102] John E. Dunn: “The UK’s 15 Most Infamous Data Breaches,” *techworld.com*, November 18, 2016.

[103] Cory Scott: “Data is not toxic - which implies no benefit - but rather hazardous material, where we must balance need vs. want,” *twitter.com*, March 6, 2016.

[104] Bruce Schneier: “Mission Creep: When Everything Is Terrorism,” *schneier.com*, July 16, 2013.

[105] Lena Ulbricht and Maximilian von Grafenstein: “Big Data: Big Power Shifts?,” *Internet Policy Review*, volume 5, number 1, March 2016. doi:10.14763/2016.1.406.

[106] Ellen P. Goodman and Julia Powles: “Facebook and Google: Most Powerful and Secretive Empires We’ve Ever Known,” *theguardian.com*, September 28, 2016.

[107] Directive 95/46/EC on the protection of individuals with regard to the processing of personal data and on the free movement of such data, Official Journal of the European Communities No. L 281/31, *eur-lex.europa.eu*, November 1995.





- [108] Brendan Van Alsenoy: “Regulating Data Protection: The Allocation of Responsibility and Risk Among Actors Involved in Personal Data Processing,” Thesis, KU Leuven Centre for IT and IP Law, August 2016.
- [109] Michiel Rhoen: “Beyond Consent: Improving Data Protection Through Consumer Protection Law,” *Internet Policy Review*, volume 5, number 1, March 2016. doi: 10.14763/2016.1.404.
- [110] Jessica Leber: “Your Data Footprint Is Affecting Your Life in Ways You Can’t Even Imagine,” *fastcoexist.com*, March 15, 2016.
- [111] Maciej Cegłowski: “Haunted by Data,” *idlewords.com*, October 2015.
- [112] Sam Thielman: “You Are Not What You Read: Librarians Purge User Data to Protect Privacy,” *theguardian.com*, January 13, 2016.
- [113] Conor Friedersdorf: “Edward Snowden’s Other Motive for Leaking,” *theatlantic.com*, May 13, 2014.
- [114] Phillip Rogaway: “The Moral Character of Cryptographic Work,” Cryptology ePrint 2015/1162, December 2015.







# 术语表



注意，本术语表的定义力求简单明了，旨在传达核心思想，可能并不完整或严谨。更多详细信息，请参阅正文中相关参考资料。

## asynchronous（异步）

不等待某件事完成（例如通过网络发送数据到另一个节点），也没有多久必须完成的事先假定。详见第5章和第8章。

## atomic（原子性）

1. 在并发操作的上下文中，描述操作在某个时间点生效，因而另一个并发进程永远不会观测到中间状态。参见隔离性。
2. 在事务上下文中，指把一组写操作逻辑绑定在一起，即使发生了故障，也保证要么全部完成，要么全部回滚。详见第7章和第8章。

## backpressure（背压）

由于接收者无法跟上，强制数据的发送者放慢发送数据速度，也称为流量控制。详见第12章。

## batch process（批处理）

一种计算过程，输入固定的（通常是很大的）数据集，并输出处理结果，但并不修

改原始输入。详见第10章。

## bounded（有界）

存在一定的上限，用以描述网络延迟（第7章）或者数据量（第11章）。

## byzantine fault（拜占庭错误）

指节点可以发生任何类型的错误，包括向其他节点发送矛盾的或者虚假消息。详见第7章。

## cache（缓存）

某个组件记住最近的数据，以加快将来可能的读取操作。通常并非数据全集，因此当缓存未命中时，必须从下层、较慢的完整数据副本或者存储系统中获取数据。

## CAP theorem（CAP理论）

一个存在很多争议的理论，目前在实践中帮助不大。详见第9章。

## causality（因果关系）

一件事发生在另一件事情之前因此构成依赖关系。例如，后面的事件是对前面（前序）事件的回复，或是建立在前序事





件之上，或者应该根据前序事件来理解。  
详见第5章和第9章。

### consensus (共识)

分布式计算的基本问题，涉及到多个节点之间达成一致（例如，哪个节点应该是主节点）。共识问题远比表面看起来要困难得多。详见第9章。

### data warehouse (数据仓库)

一个数据库，其中来自多个不同的OLTP系统的数据已被合并，主要用于分析目的。详见第3章。

### declarative (声明式)

描述一些东西应有的属性，而不描述关于如何实现的确切步骤。在查询上下文中，查询优化器采用声明式查询并决定如何最优执行。详见第2章。

### denormalize (反规范化)

通常以缓存或索引的形式引入了一定冗余或重复数据集，用来加速读取。一个非规范化值是一种预先计算好的查询结果，类似于实体化/物化视图。详见第7章和第11章。

### derived data (派生数据)

由其他数据通过可重复的过程所创建出的数据集，如有必要可以多次运行。通常用来加速对特定类型数据的读取访问。索引、缓存和物化视图属于典型的派生数据。详见第三部分。

### deterministic (确定性)

给定相同的输入总是得到相同的输出。这意味着它不能依赖于任何随机数字、时间、网络通信或其他不可预知的事情。

### distributed (分布式)

通过网络连接起来的多节点。以部分失效为主要特征：系统的某些部分可能会发生错误，但其他部分仍然正常工作，软件通常无法确切知道究竟发生了是什么样的错误。详见第8章。

### durable (持久化)

以某种方式存储数据，即使发生一些错误，也通常认为不会丢失数据。详见第7章。

### ETL

即提取-转换-加载。从源数据库提取数据，将其转换为更适合分析查询的格式，并加载进数据仓库或批处理系统。详见第3章。

### failover (切换)

对于单个主节点系统，将主节点角色从一个节点转移到另一个节点的过程。详见第5章。

### fault-tolerant (容错)

如果组件发生错误（例如节点崩溃、网络中断等），可以自动从中恢复。详见第1章。

### flow control (流控)

参见backpressure。

### follower (从副本)

该副本不直接接受客户端的写请求，而是从主节点（主副本）上接收数据变化。详见第5章。

### full-text search (全文搜索)

任意关键字搜索，通常支持相似度匹配或者通配符等。全文索引是某种二级索引用来支持这种查询。详见第3章。

### graph (图)

一种数据结构，包含顶点（节点/实体对象）和边（描述对象间关系）。详见第2章。

### hash (哈希)

一个函数用于将输入转为看似随机的输出。相同的输入一定产生相同的输出；不同的输入通常对应不同的输出，尽管有可能产生相同的输出（即冲突）。详见第6章。

### idempotent (幂等)

描述操作可以安全地重试，多次执行的结



果和单次执行完全相同。详见第11章。

### index (索引)

一种加速查询的数据结构，输入特定的值查询所有符合条件的记录。详见第3章。

### isolation (隔离性)

在事务上下文中指并发执行的事务可能发生交互的程度。可串行化隔离是最强的保证，但弱隔离也广泛使用。详见第7章。

### join (联结)

把某种关联的数据联结在一起，常用于某条记录又引用另外一条（外键，文档引用，图中的边等），查询时需要同时得到所指向的记录。详见第2和第10章。

### leader (主节点，主副本)

当数据或者服务在多个节点间复制时，指定主节点来负责数据更改。主节点可以通过某种协议选举，或者由管理员手工指定。详见第5章。

### linearizable (线性化)

对外表现得好像好像只有单一副本，且更新都以原子方式进行。详见第9章。

### locality (局部性)

性能优化的手段，如果数据同时频繁地访问，则把一些数据放在相邻的位置。详见第2章。

### lock (锁)

一种机制来确保某个对象只能被一个线程、节点、事务所访问；如果还有其他待访问者，则必须等锁释放之后才能继续。详见第7章和第9章。

### log (日志)

一种追加修改的文件用来保存数据。预写日志用来帮助存储引擎在崩溃之后进行恢复（第3章）；日志式存储则以追加日志作为主存方式（第3章）；复制日志用来把更新从主节点拷贝到从节点（第5章）；事件日志用来表示一个持续的事件流（第11章）。

### materialize (实体化，物化)

把数据预先计算并保存结果，而不是请求到来时按需计算。详见第3章和第11章。

### node (节点)

软件运行所在的计算机实例，可以通过网络与其他节点一起完成某些工作。

### normalized (规范化)

一种结构范式，其中没有冗余或者重复的数据；对于规范化的数据库，更改时只需修改一处，而没有多个地方的多个副本需要修改。详见第2章。

### OLAP (联机分析处理)

联机分析型处理。常见的访问模式包括对大量数据进行聚合（计数，求和，平均值等）操作，详见第3章。

### OLTP (联机事务处理)

联机事务型处理。常见访问模式主要是快速查询，或者修改特定的小部分数据，多通过主键来索引。详见第3章。

### partitioning (分区)

把那些无法由单台机器完成的超大数据集或者计算任务分解为更小的部分，然后分布到多台机器上，又称为分片。详见第6章。

### percentile (百分位数)

一种统计数据值分布的方法，描述有多少数据在一个给定阈值之上（或之下）。例如p95 响应时间为 $t$ 表示，一段时间内有95%的请求响应时间小于 $t$ ；而剩下的5%则高于 $t$ 。详见第1章。

### primary key (主键)

一个值（数字或者字符串）用来唯一标识一条记录。许多应用系统中，主键由系统在记录创建时自动生成（采用顺序递增方式或者随机方式），而不是由用户指定。参考二级索引。

### quorum (法定票数，仲裁)

最小有效的节点参与投票数，低于此数则投票结果无效。详见第5章。



### rebalance (再平衡)

在多个节点之间移动数据或服务,使得负载更加均衡。详见第6章。

### replication (复制)

把相同的数据存放在多个节点上,这样即使一些节点出现问题,其他节点上的副本可以继续访问。详见第5章。

### schema (模式)

一种描述数据的结构,通常包括元素组成,数据类型等。可以动态监测数据是否符合模式要求(详见第3章),模式有可能会发生变化(详见第4章)。

### secondary index (二级索引)

在主数据之外所维护的辅助数据,可以有效帮助查询符合特定条件的记录。详见第3章和第6章。

### serializable (可串行化)

一种保证,即多事务并发执行的结果与单个串行执行的结果相同。详见第7章。

### shared-nothing (无共享架构)

一种系统架构包括多个节点,每个节点各有独立的CPU、内存和磁盘,节点通过常用网络连接。与之对应的是共享内存架构和共享磁盘架构。详见第二部分。

### skew (倾斜)

1. 多个分区间出现负载不均衡,因而某些分区承担了太多的数据或者请求,也称为热点。详见第6章和第10章。
2. 一种时序异常,导致事件呈现非预期的,违背顺序关系。详见第7章的读倾斜、写倾斜和时钟倾斜。

### split brain (脑裂)

同时有两个节点都自认为是系统主节点的异常情况。会违背系统正确性保证。详见第5章和第8章。

### stored procedure (存储过程)

一种将事务逻辑打包的方式,这样可以在服务器端高效执行,而不用与客户端出现反复的网络监护。详见第7章。

### stream process (流式处理)

一种持续处理无边界数据事件的计算方式。详见第11章。

### synchronous (同步)

异步的反义词。

### system of record (记录系统)

拥有原始、权威数据的系统,也被称为source of truth。可以对其修改,然后更多的派生数据可以基于此而继续处理。详见第三部分。

### timeout (超时)

一种检测失效的最简单的方式。即在规定时间内观察节点的响应。但无法区分究竟是节点问题还是网络问题。详见第8章。

### total order (全序关系)

一种总是满足事物可比较的方法;如果只有一部分满足可比较,则称之为偏序。详见第9章。

### transaction (事务)

将多个读写操作聚合为一个逻辑操作,以简化错误处理和并发性。详见第7章。

### two-phase commit (两阶段提交)

一种算法用来协调多个节点完成事务,要么全部提交,要么全部中止。详见第9章。

### two-phase locking (两阶段加锁)

一种算法用以实现串行化隔离,事务将获取所有读/写对象的锁,直到事务完成结束。详见第7章。

### unbounded (无界)

有界(bounded)的反义词,不存在可知的上限。

## 作者介绍

---

**Martin Kleppmann**是英国剑桥大学分布式系统方向的研究员。此前，他曾是LinkedIn和Rapportive等互联网公司的软件工程师，负责大规模数据基础设施建设。在此过程中他遇到过一些困难，因此他希望这本书能够帮助读者避免重蹈覆辙。

Martin还是一位活跃的会议演讲者、博主和开源贡献者。他认为，每个人都应该学习深刻的技术理念，对技术的深入理解能帮助我们开发出更好的软件。

## 封面介绍

---

本书封面上的动物是一头印度野猪，它是在印度、缅甸、尼泊尔、斯里兰卡和泰国发现的一种野猪的亚种。它们与欧洲的公猪不同，有较高的后髯，没有毛绒底毛，还有一个更大、更直的头骨。

印度野猪有一头灰色或黑色的头发，脊椎上有僵硬的硬毛。雄性有突出的犬齿（或称獠牙），用来与对手战斗或抵御掠食者。雄性要比雌性格格大，但是整体平均肩高为33~35英寸，体重200~300磅。它们的天敌包括熊、老虎和各种大型猫科动物。

这些动物属于夜行杂食性动物，它们吃各种各样的东西，包括树根、昆虫、腐肉、坚果、浆果和小动物。野猪也以翻拱垃圾和农田为大家所熟知，造成大量破坏，而招致农场主的憎恨。它们每天需要食用4000~4500卡路里的热量。公猪拥有发达的嗅觉，能够帮助它们寻找地下植物和穴居动物。但是，它们的视力很差。

野猪在人类文化中一直具有重要意义。在印度教传说中，野猪是毗湿奴神的化身。在古希腊遗址中，它是一个勇敢的失败者的象征（与胜利的狮子相反）。由于它的侵略性，它被描绘在斯堪的纳维亚、日耳曼和盎格鲁-撒克逊战士的盔甲和武器上。在中国十二生肖中，它显得比较急躁。

O'Reilly出版物封面上的许多动物都属于濒危动物；所有这些动物对世界来说都是重要的。如果你想了解更多关于如何为它们提供帮助，请访问[animals.oreilly.com](http://animals.oreilly.com)。

本书封面图片来自Shaw的*Zoology*。



# 数据密集型应用系统设计

当今很多应用系统的核心挑战在于数据。有诸多系统层面的问题亟须解决方案，例如扩展性、一致性、可靠性、效率问题，以及可维护性等。此外，摆在我们面前的还有各种各样的工具，包括关系型数据库、NoSQL存储系统、批处理、流处理，以及各种消息系统。问题是，哪些是最合适的工具？又该如何评估、考量这些林林总总商业名词背后的技术呢？

本书内容详实并富有实际指导价值，本书作者将带你纵览各种数据处理和存储技术，探讨背后的优劣与取舍之道。软件应用虽有千般变化，却终有若干原则贯穿其中。通过本书，软件开发者和架构师将学到如何将这些原则用于实践，以及如何在最新的应用架构中充分发挥数据的威力。

- 深入探索常用分布式系统内部机制，学习如何高效运用这些技术。
- 分析各种工具的优势和不足，帮助做出明智的设计决策。
- 解析一致性、扩展性、容错和复杂度之间的权衡利弊。
- 介绍分布式系统研究的最新进展（现代数据库的基石）。
- 揭示主流在线服务的基本架构。

Martin Kleppmann是英国剑桥大学分布式系统方向的研究员。此前，他曾是LinkedIn和Rapportive等互联网公司的软件工程师，负责大规模数据基础设施建设。Martin还是一位活跃的会议演讲者、博主和开源贡献者。

“这本书非常棒！它弥补了分布式理论和工程实践之间的鸿沟。我多么希望十年前它就能面市，这样我可以早点读到，从而避免这一路犯下的许多错误。”

——Jay Kreps

Apache Kafka的创造者，  
Confluent CEO

“本书应该是软件开发者的必读书籍。很少有这样的书籍可以如此完美地衔接理论与实践，对于那些设计和实现基础数据架构的开发者而言，它可以帮你做出睿智的决定。”

——Kevin Scott

微软CTO

DATA | HADOOP

O'Reilly Media, Inc. 授权中国电力出版社出版

此简体中文版仅限于在中华人民共和国境内（但不允许在中国香港、澳门特别行政区和中国台湾地区）销售发行  
This Authorized Edition for sale only in the territory of People's Republic of China (excluding Hong Kong, Macao and Taiwan)



ISBN 978-7-5198-2196-8



定价：128.00元